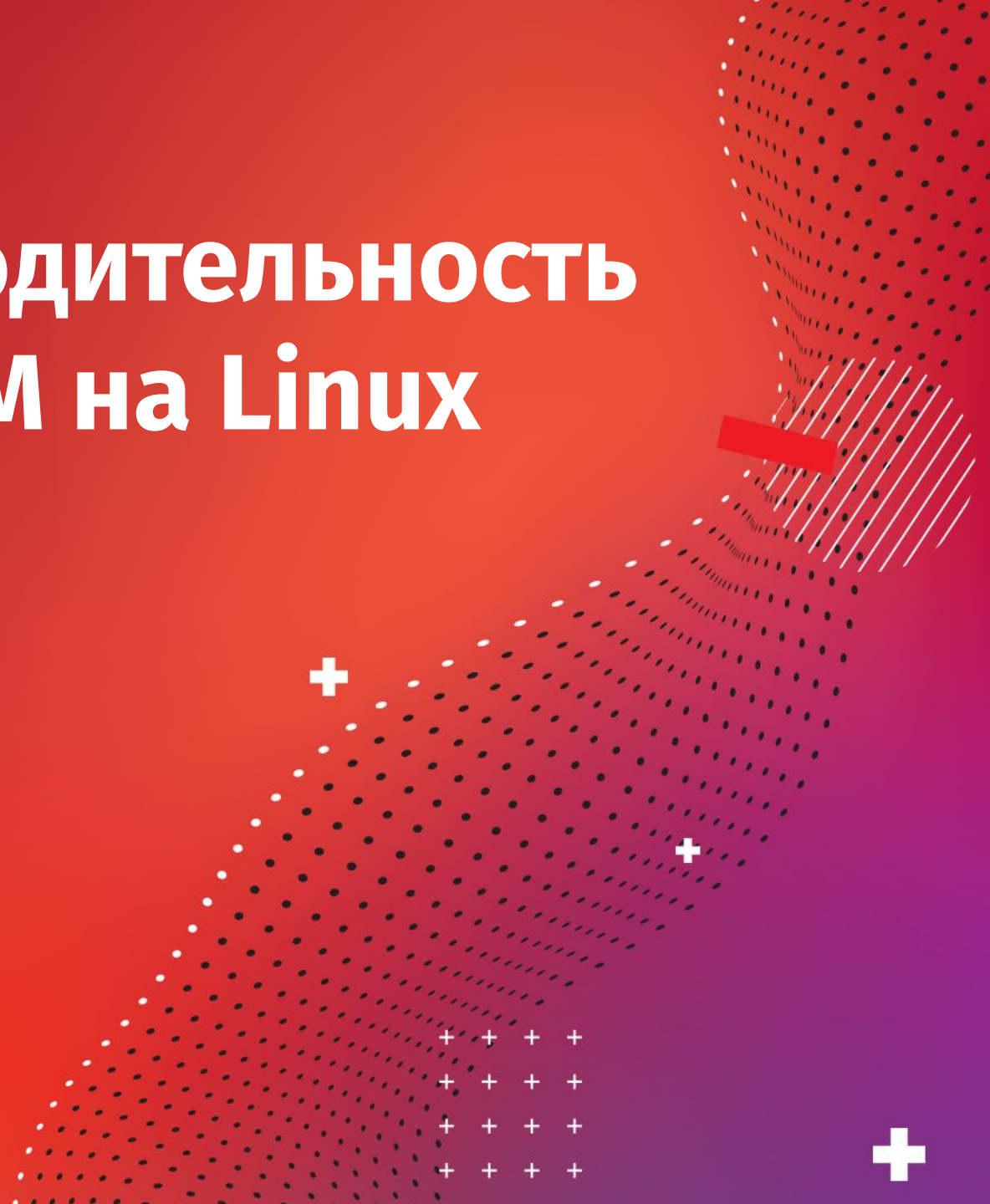


Повышаем производительность файлового I/O в JVM на Linux

Дмитрий Бундин



HighLoad++
Весна 2021



О спикере

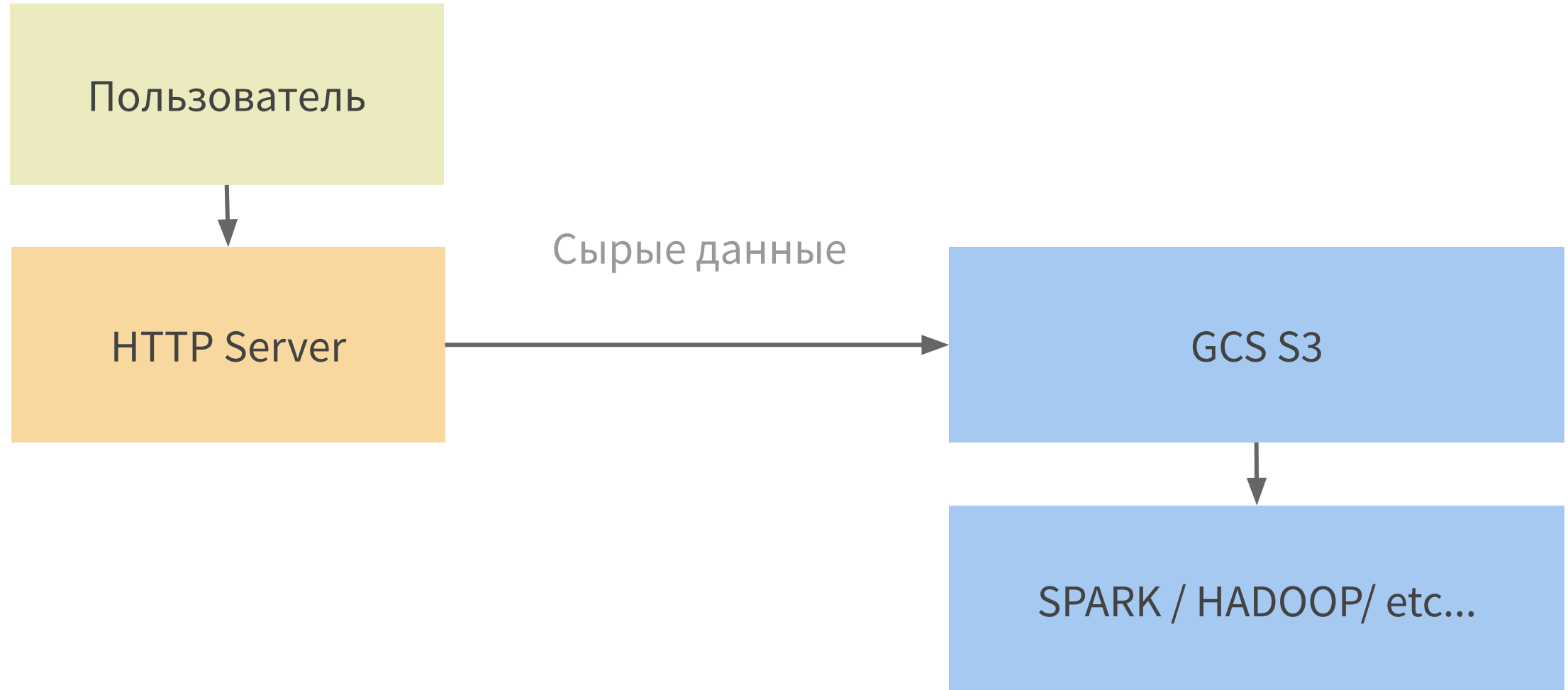


Дмитрий Бундин

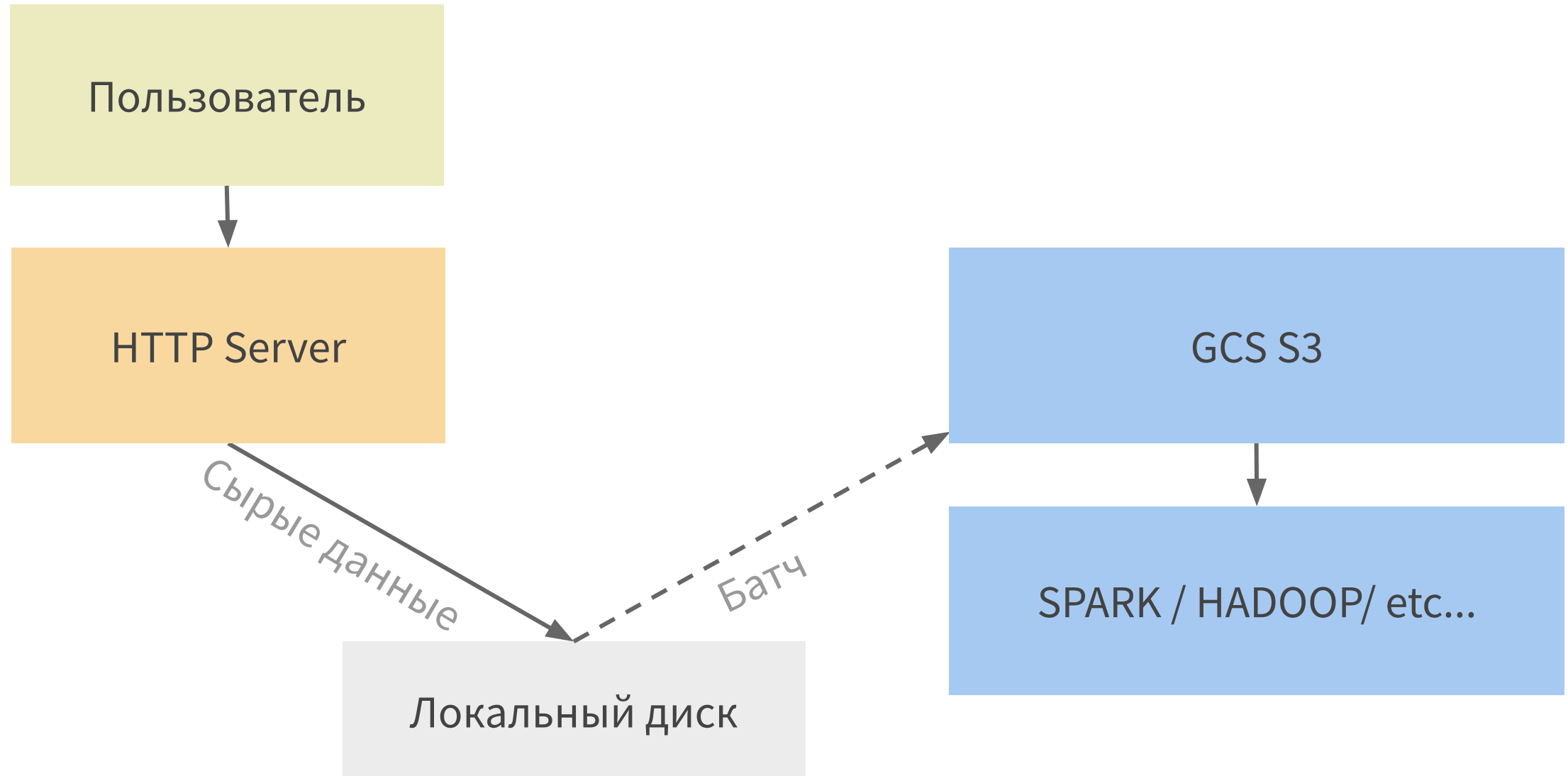
Senior Big Data developer

Дмитрий в ИТ с 2014 года. В последние несколько лет активно занимается разработкой I/O интенсивных приложений на Java, C/Linux и связанными с ними вопросами производительности. До этого разрабатывал банковский софт и платформу обработки данных на Spark в сфере рекламы. Имеет опыт работы со Scala, функциональным программированием и typelevel-стеком. В настоящее время является старшим Big Data-разработчиком в Grid Dynamics.

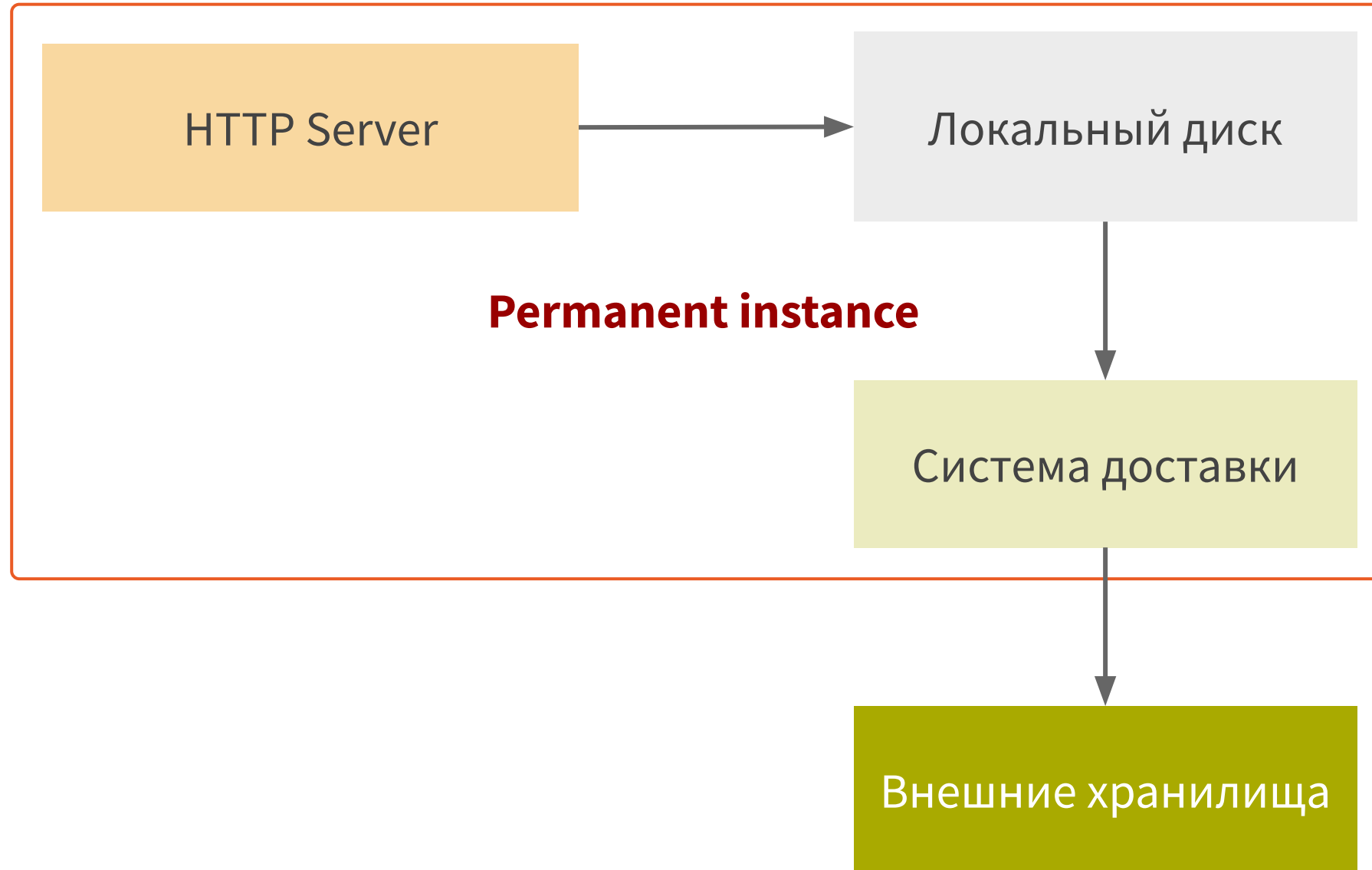
Проект в сфере рекламы



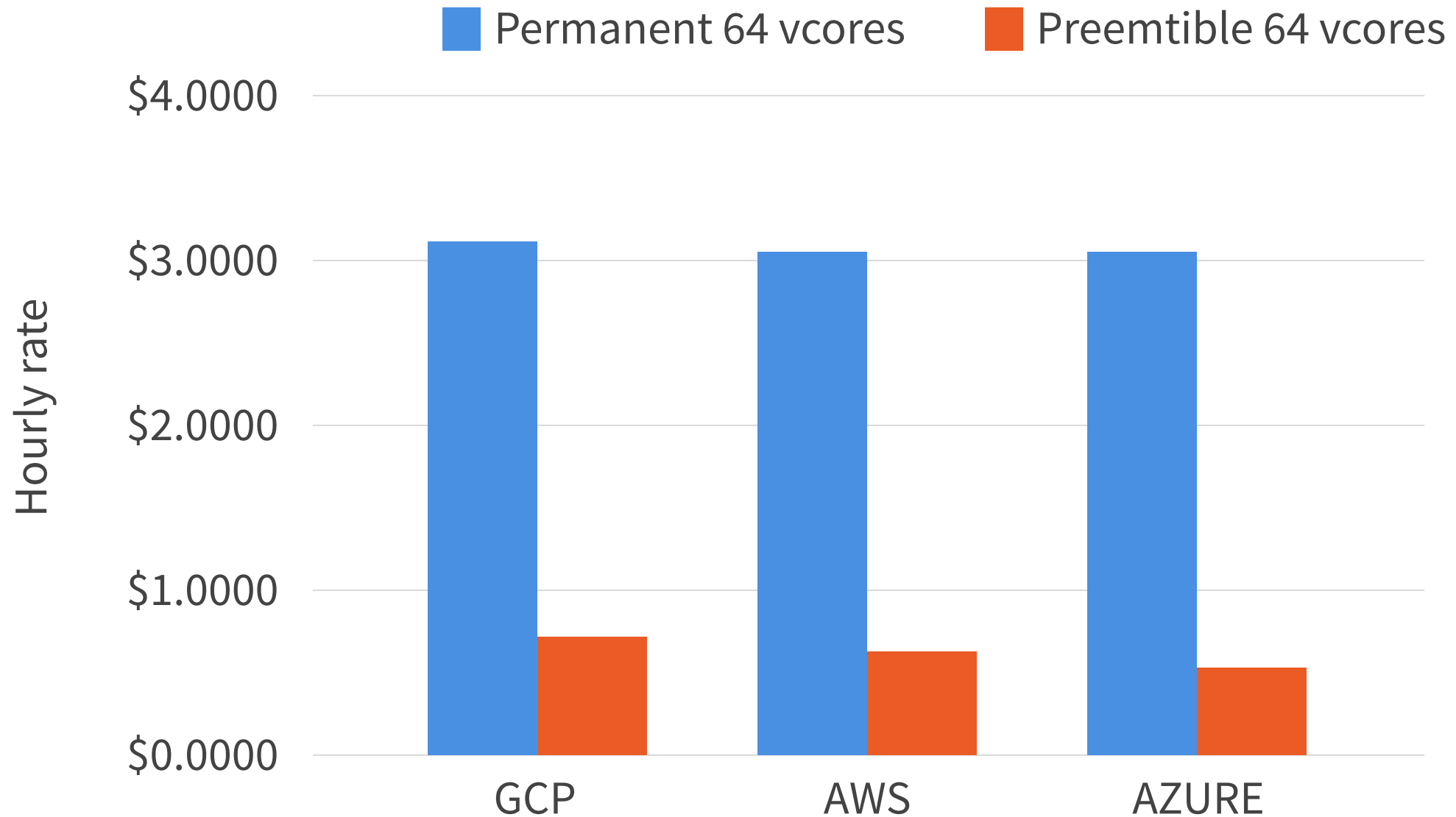
Проект в сфере рекламы



Проект в сфере рекламы: выключение инстанса



Permanent vs Preemptible: стоимость



Снижение затрат на инфраструктуру

Permanent Instances

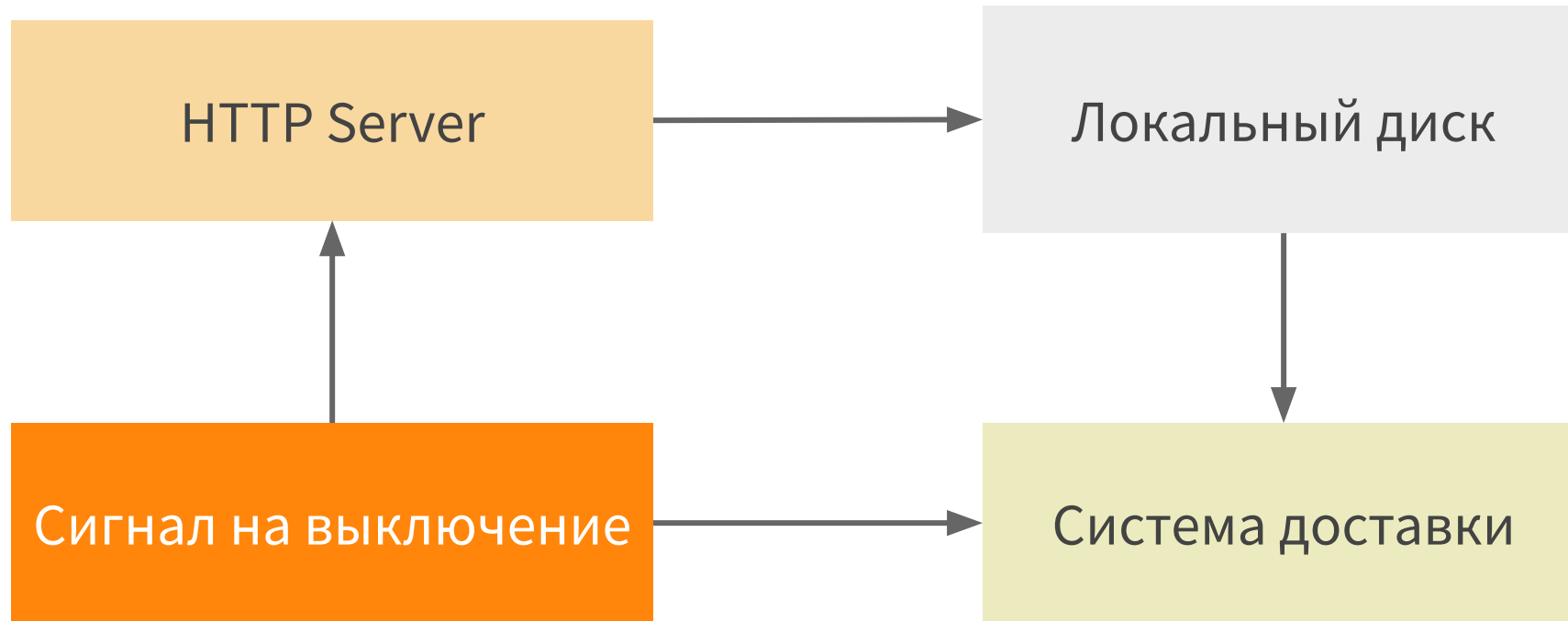


Высокая стоимость

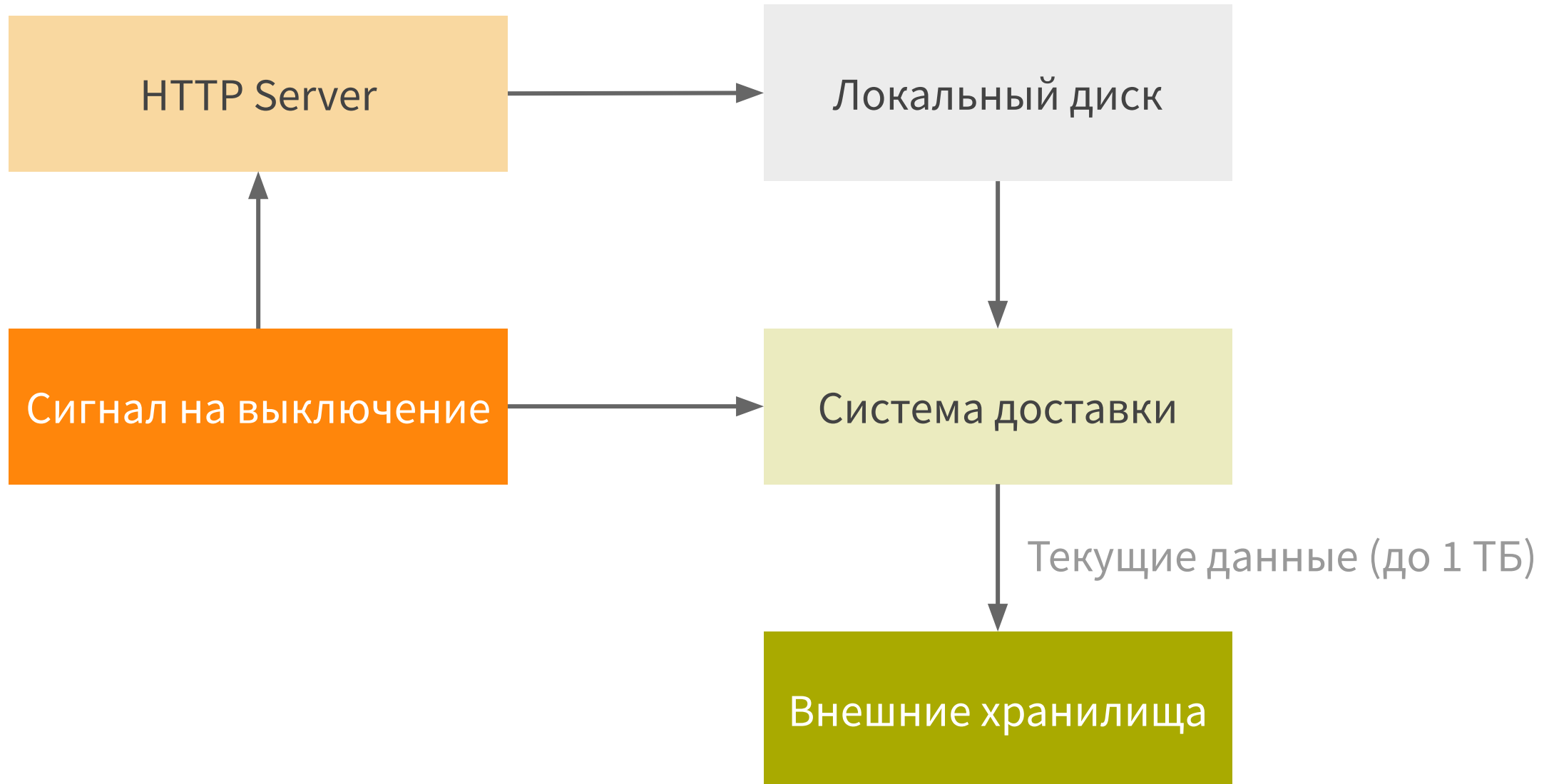


Spot Instances

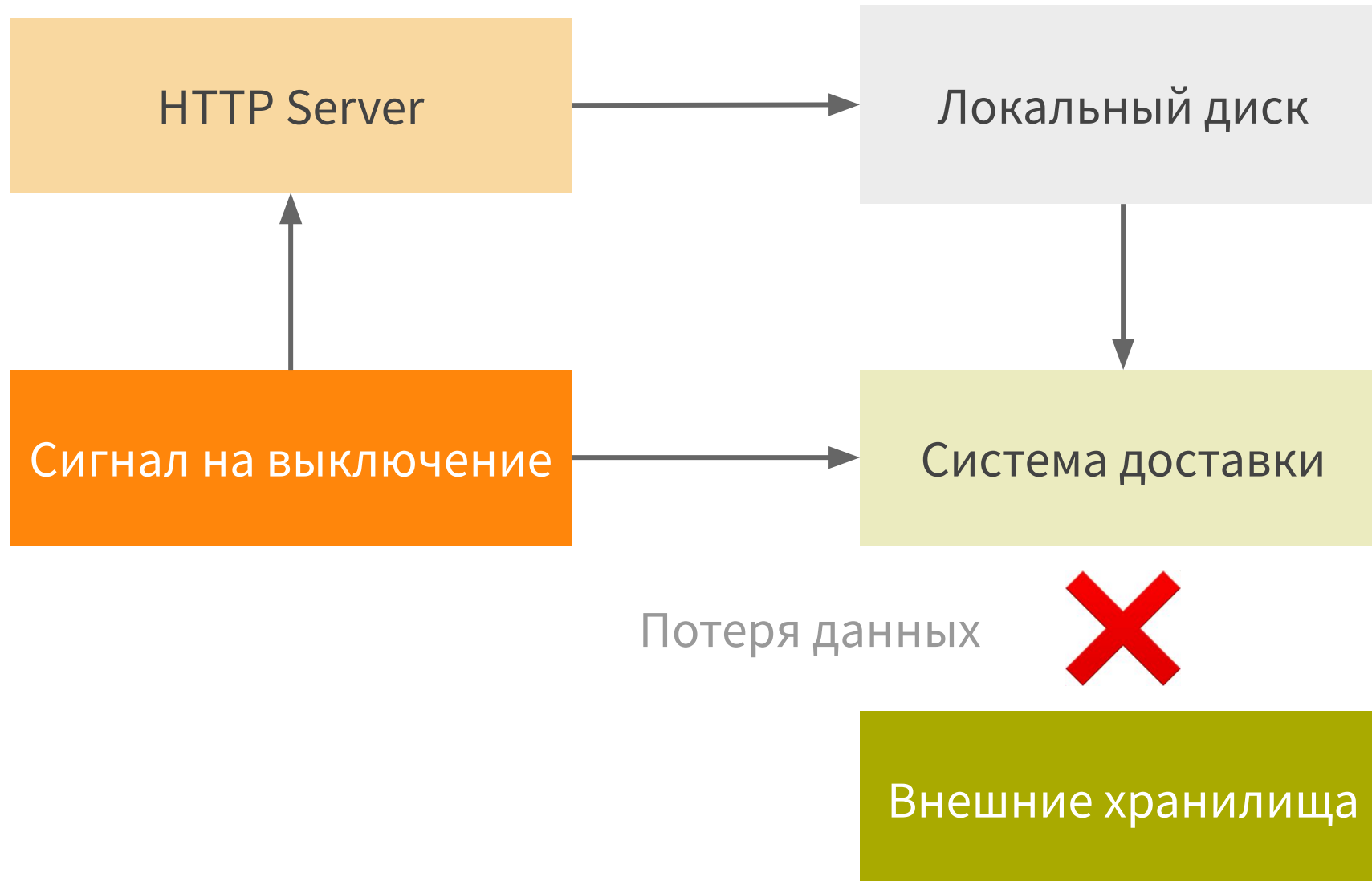
Проект в сфере рекламы: выключение инстанса



Проект в сфере рекламы: выключение инстанса



Проект в сфере рекламы: выключение инстанса



Задачи проекта

Доставка данных “сервис” → “внешние хранилища”

```
graph TD; A[Доставка данных “сервис” → “внешние хранилища”] --> B[Уменьшить вероятности потери]; A --> C[Оптимизировать доставку];
```

Уменьшить
вероятности потери

Оптимизировать
доставку

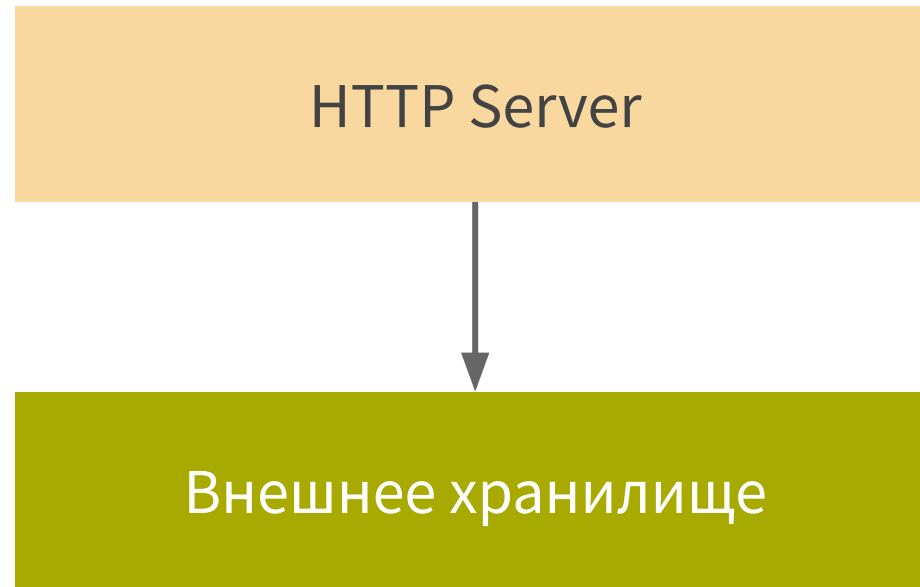
Уменьшение вероятности потери данных

1 1 0 1 1 1 0 1 1 1 1
0 1 0 1 1 1 1 0 0 0 0 1 1 1 1 1 1 0 1
1 0 1 0 1 0 1 1 1 0 1 1 1 1 0 1 0 1 1 1 1 0 1
1 1 0 1 0 0 0 0 1 0 1 1 0 1 1 0 1 1 1 1 1 0
1 0 1 0 1 0 1 1 0 1 0 1 1 1 1 0 0 1 0
1 0 1 0 1 0 1 0 1 1 1 1 1 0 1 0 1 1 0 0 1 0 1
0 1 0 1 0 1 1 1 0 1 0 1 1 0 0 1 0 1 1 1 1 0
0 1 1 1 1 1 0 1 0 1



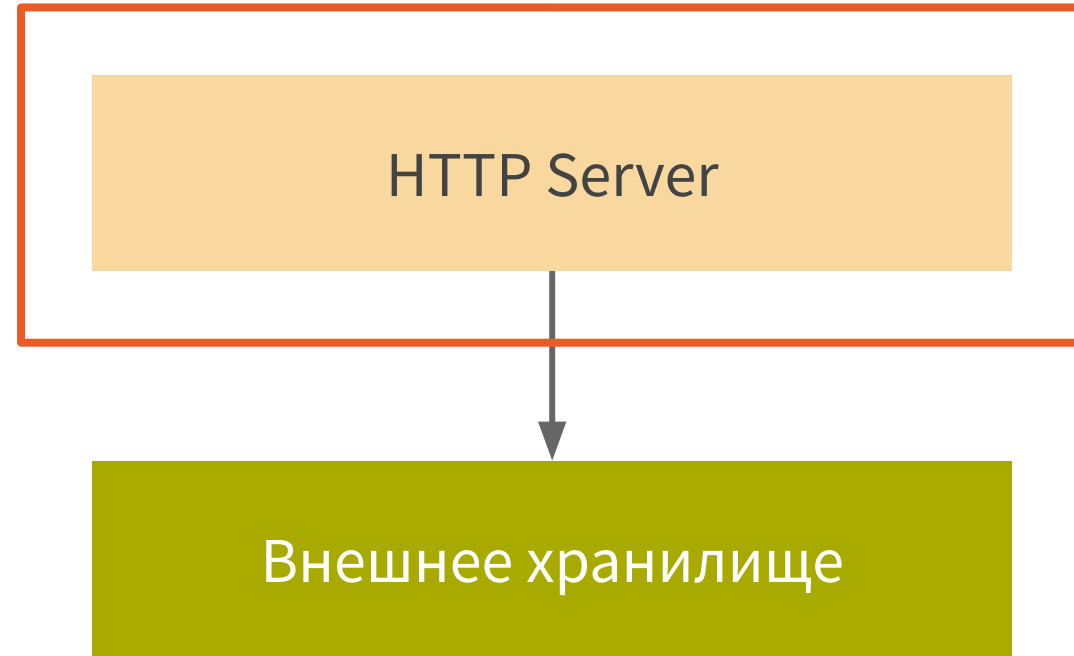
0 1 1 1 0 1 0 1 0 1
0 1 0 1 0 1 1 1 1
0 0 1 0 1 0 1
0 1 0 1 0 0 1 0 1
0 1 0 1 0 1 0 1 0
0 1 1 0 1 1 1 1 0
1 0 1 0 1 0 1 0 1 1 1
1 1 0 0 1 0 1
0 1 0 0 0 0 1 1 1 1 0

Стриминг данных напрямую



Стриминг данных напрямую

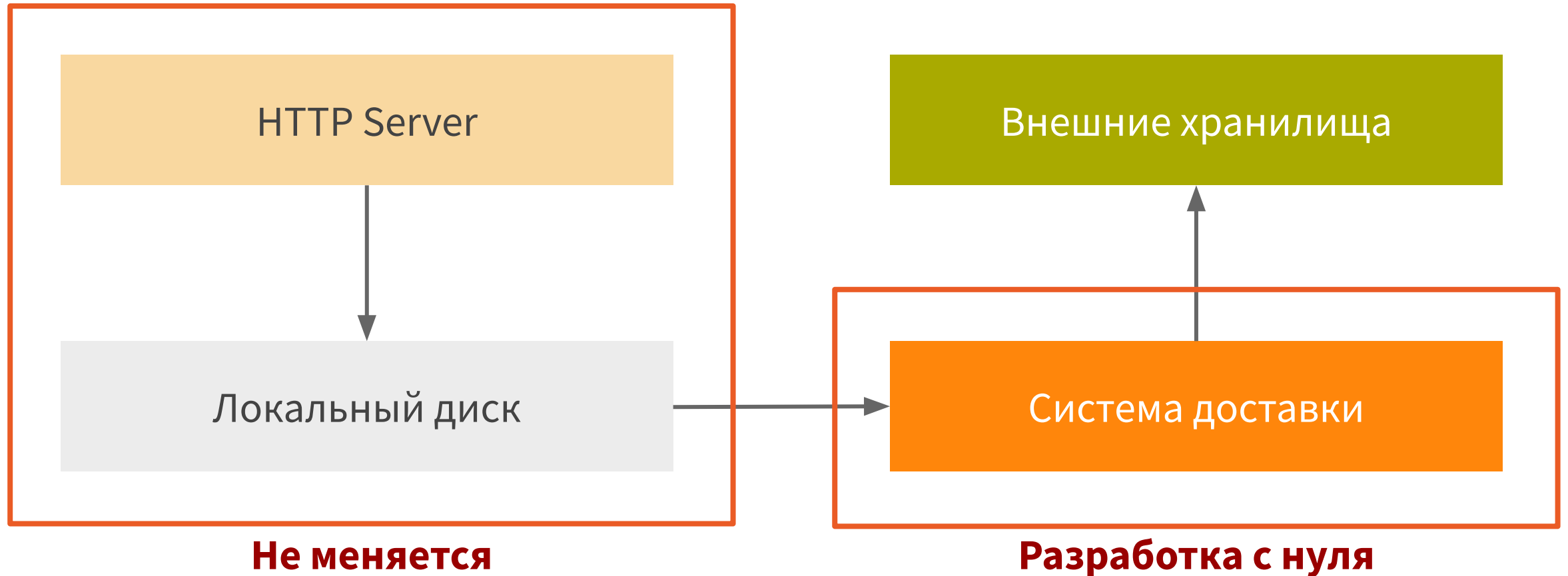
Требуется доработка



Стриминг данных с локальных дисков



Стриминг данных с локальных дисков



#2 Стриминг данных: сравнение

Стриминг данных напрямую

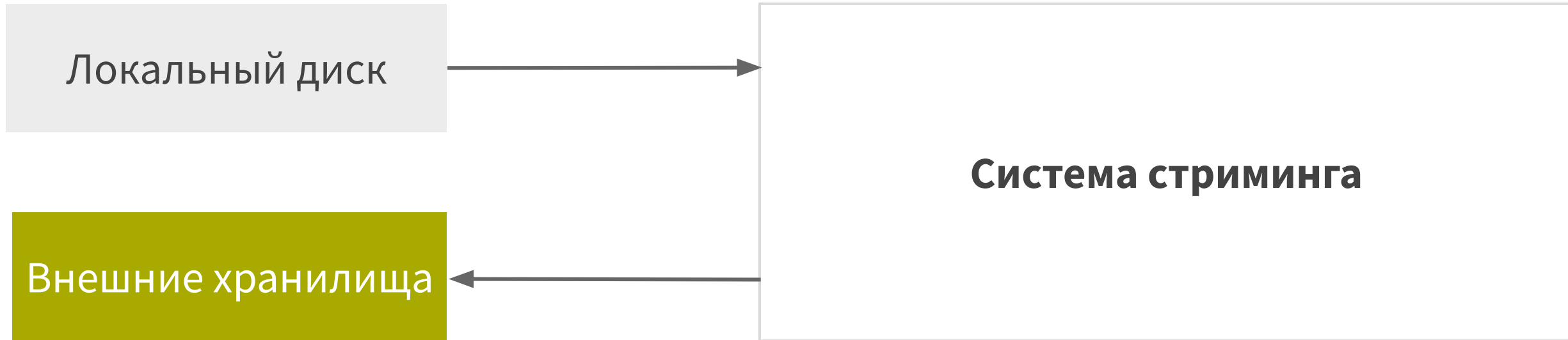
- Допиливание HTTP Server
- Вероятность потери сохраняется

Стриминг данных с дисков

+ Гибкость

- ~~- Допиливание HTTP Server~~
- ~~- Вероятность потери сохраняется~~

Схема системы стриминга



Оптимизация доставки данных

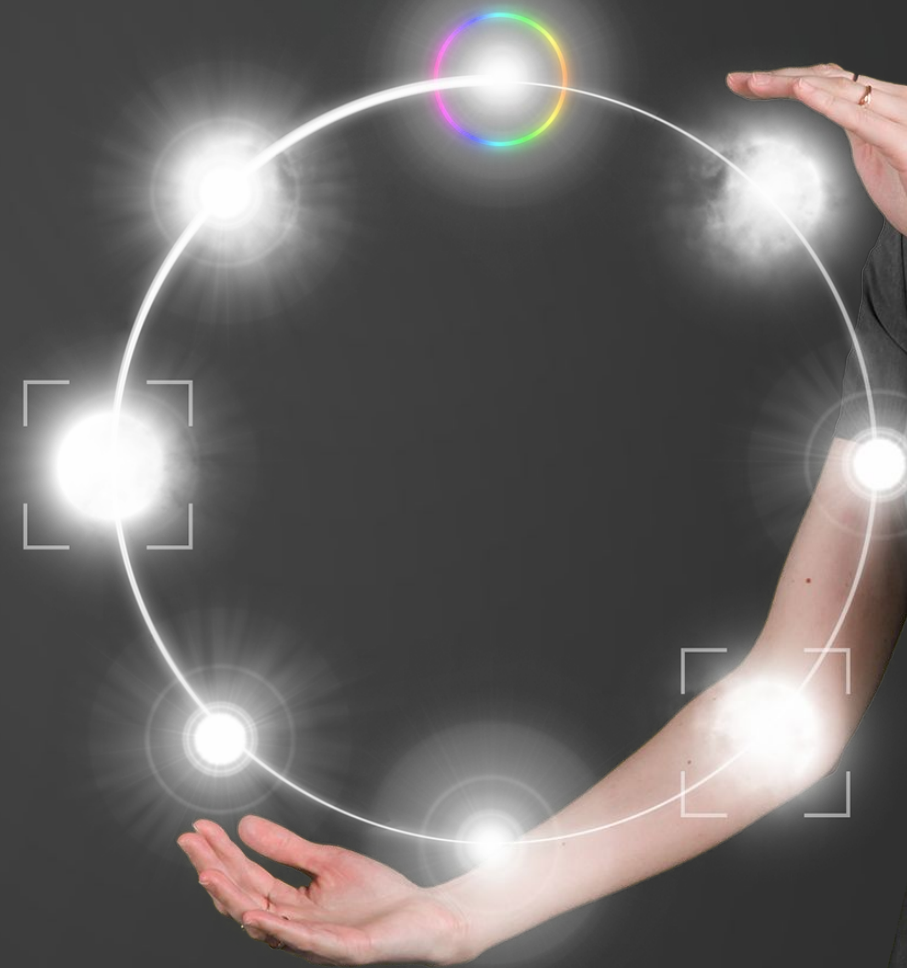


Схема системы стриминга

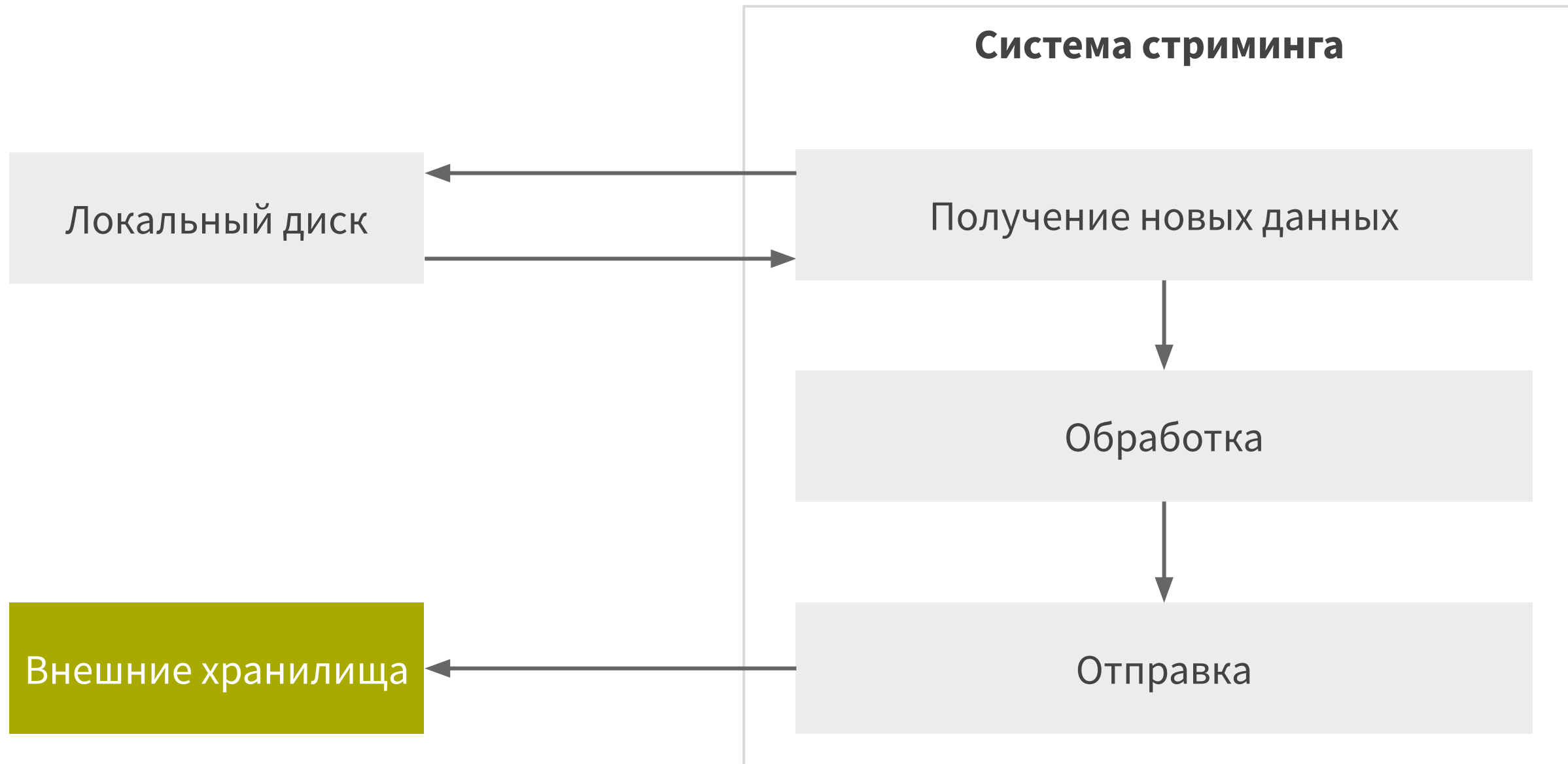


Схема системы стриминга

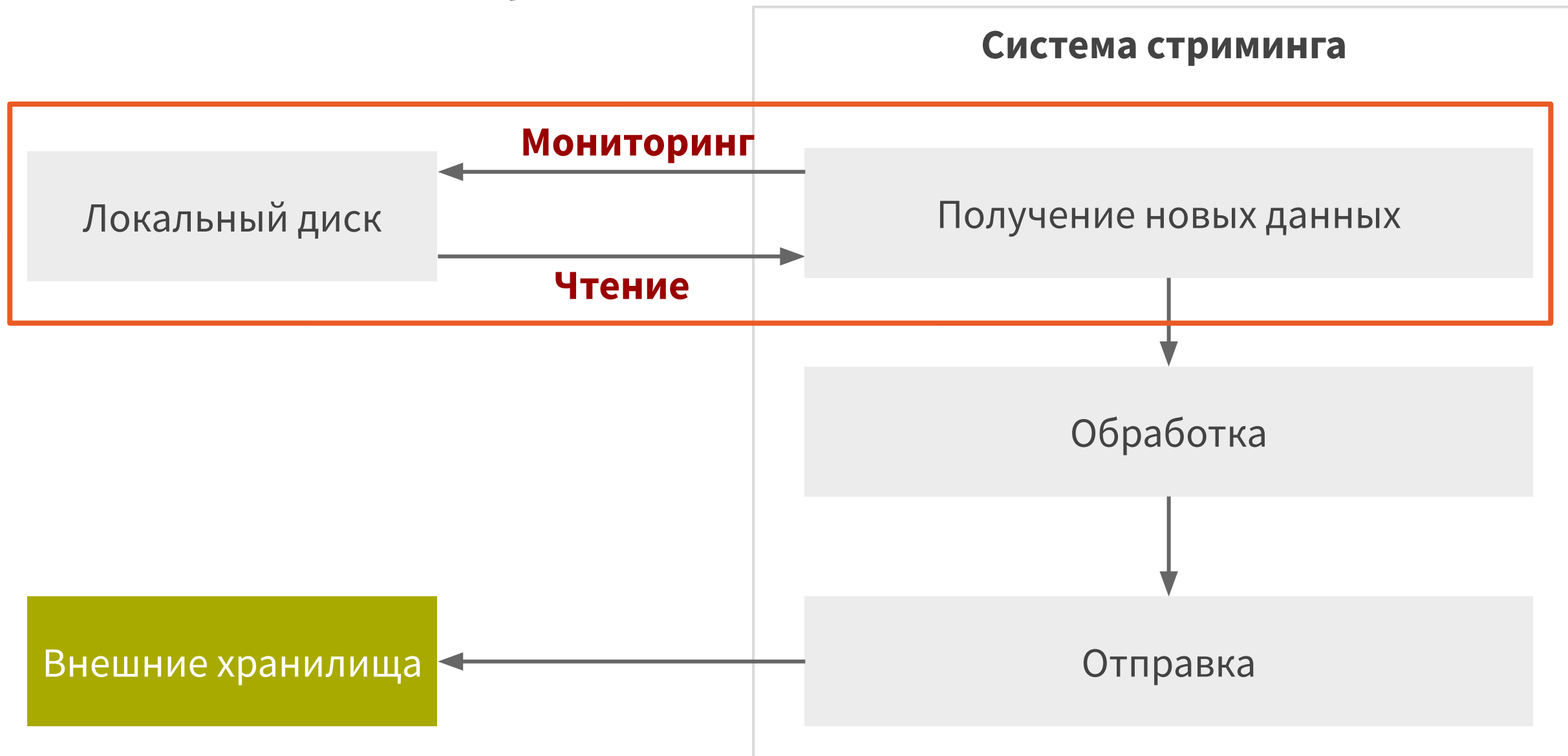
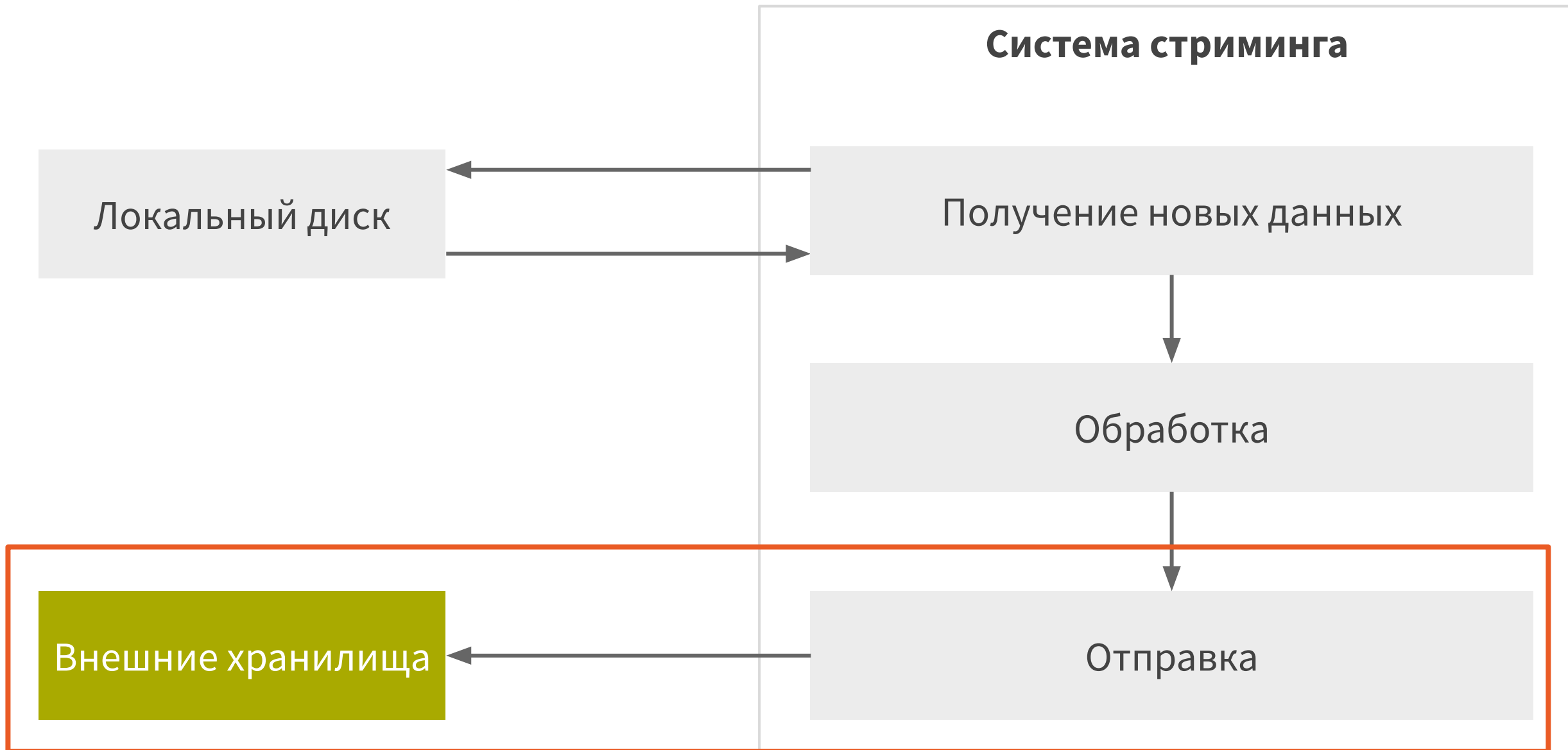


Схема системы стриминга



Схема системы стриминга



Система стриминга: Файловое I/O



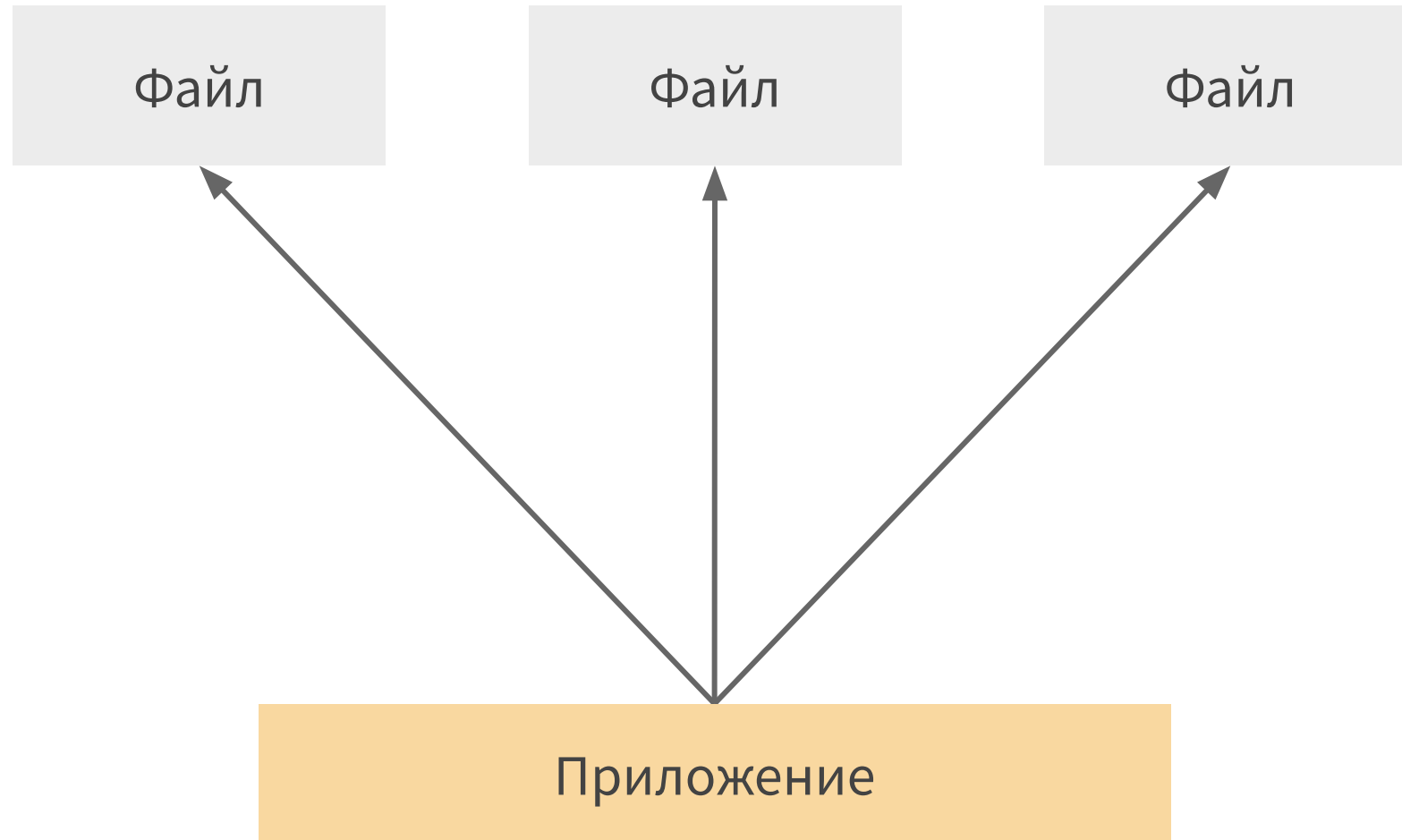
Система стриминга: Файловое I/O

- Определение изменившихся файлов
- Определение предыдущей позиции стриминга
- Чтение новых данных

Система стриминга: Файловое I/O

- **Определение изменившихся файлов**
- Определение предыдущей позиции стриминга
- Чтение новых данных

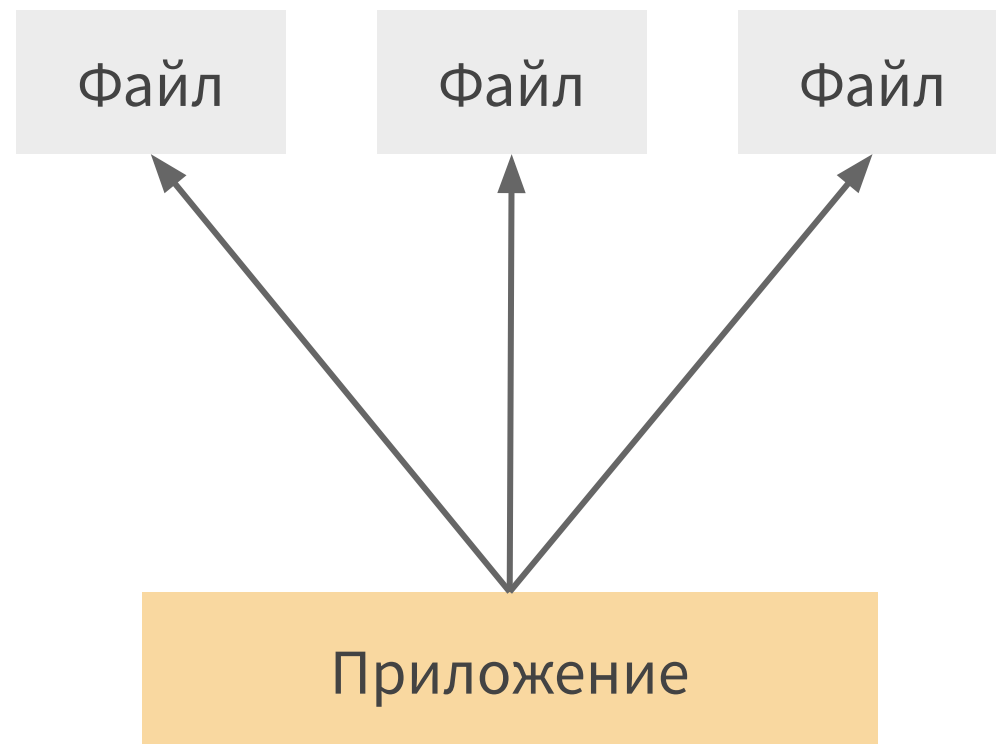
Поллинг



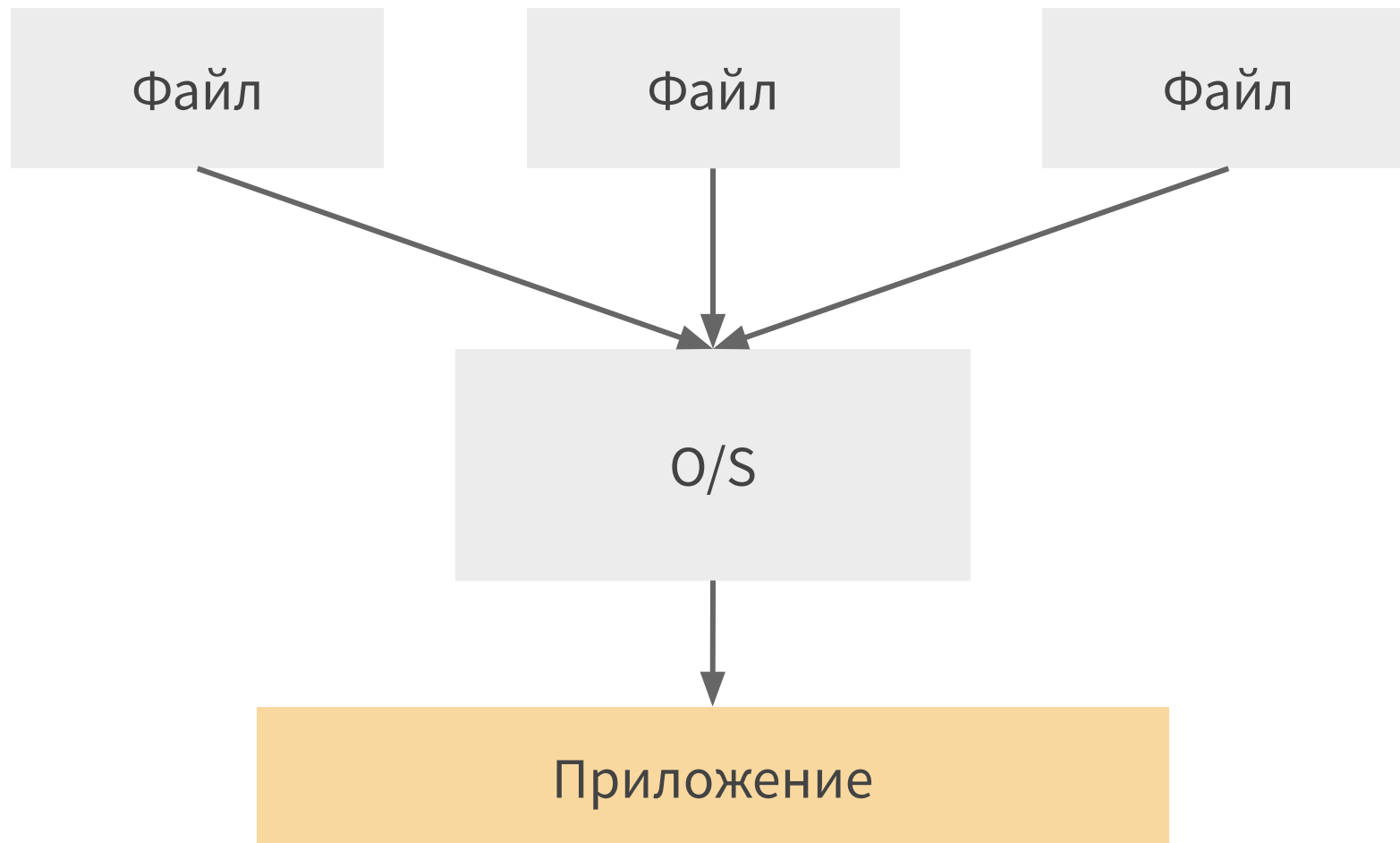
Поллинг

Недостатки:

- Poll с фиксированным интервалом
- Затраты CPU на неизменившиеся файлы



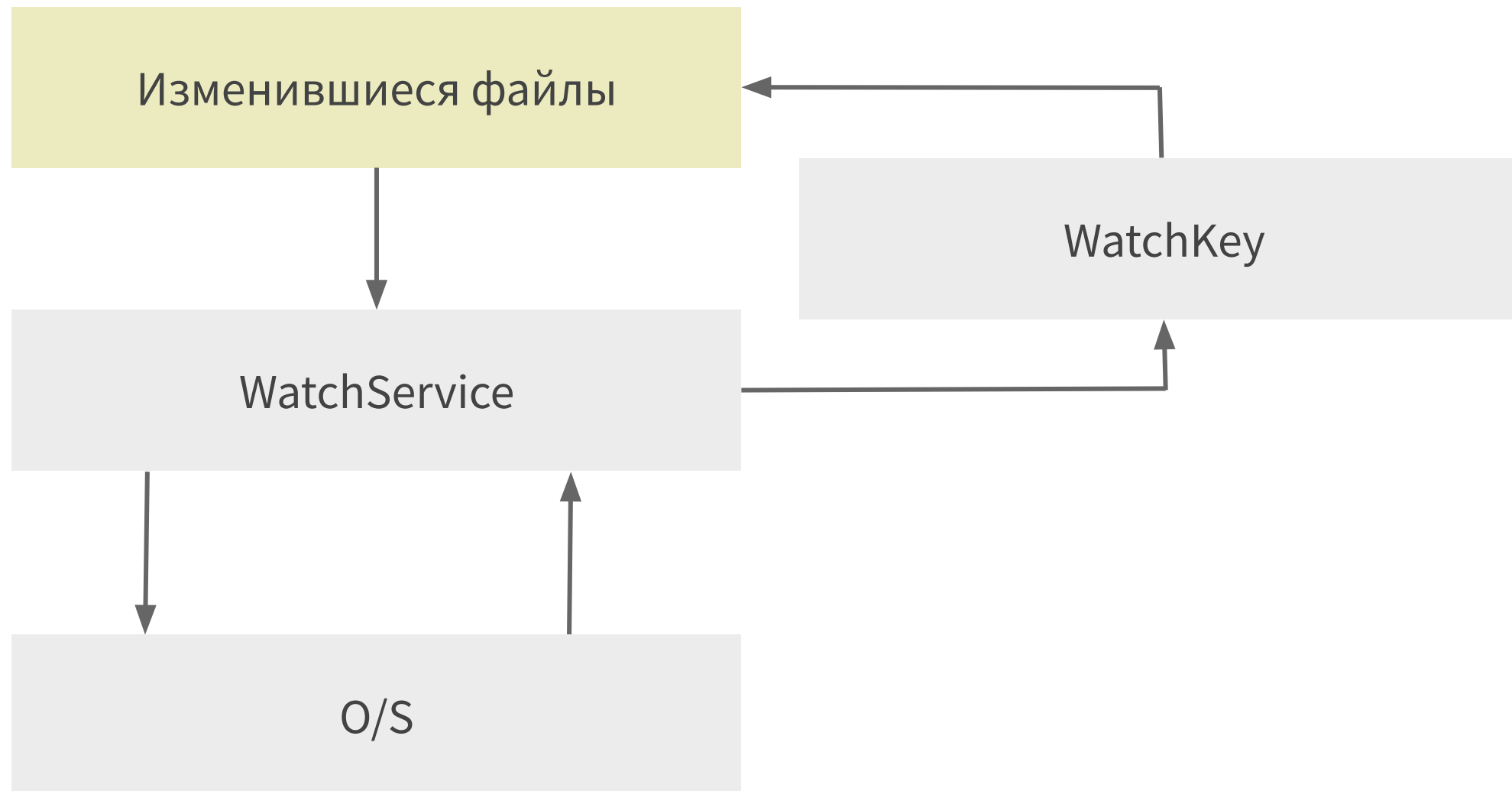
Получение уведомлений O/S



Получение уведомлений O/S: JDK java.nio.file

```
1 package java.nio.file;
2
3 public interface WatchService {
4     WatchKey poll(long timeout, TimeUnit unit)
5 }
6 ...
7 public interface WatchKey {
8     List<WatchEvent<?>> pollEvents();
9 }
10 ...
11 public interface WatchEvent<T> {
12     T context();
13 }
```

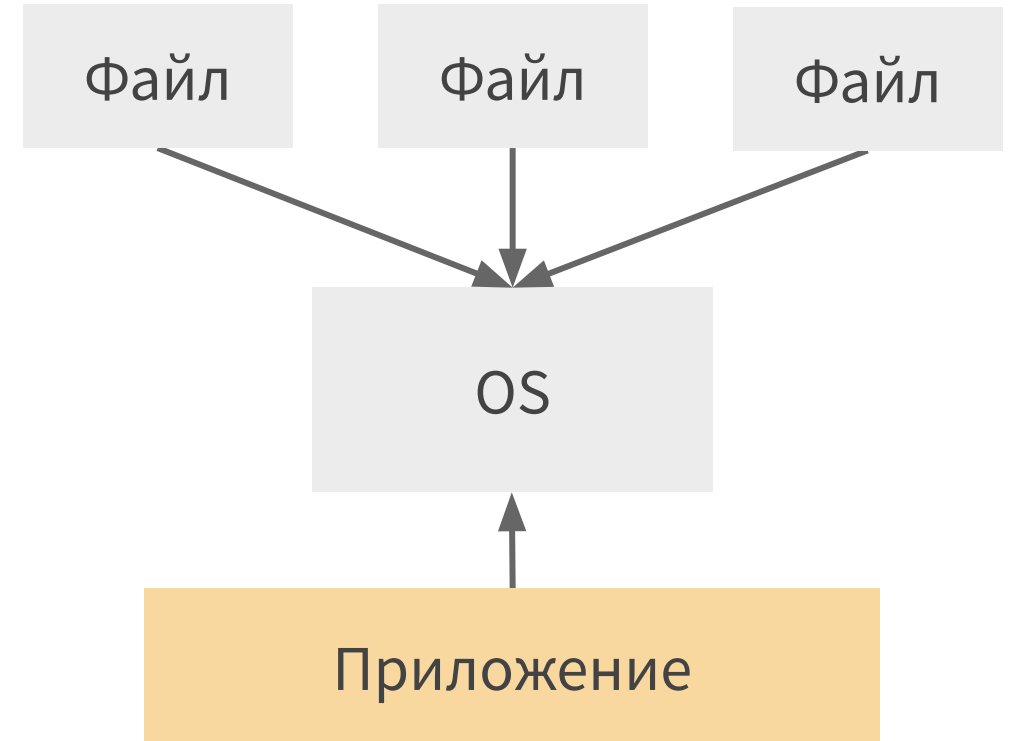
JDK java.nio.file: схема использования



Получение уведомлений O/S: JDK java.nio.file

Недостатки:

- Скучное множество доступных эвентов



JDK java.nio.file: доступные события

```
1 package java.nio.file;
2
3 public final class StandardWatchEventKinds {
4
5     public static final WatchEvent.Kind<Path> ENTRY_CREATE
6
7     public static final WatchEvent.Kind<Path> ENTRY_DELETE
8
9     public static final WatchEvent.Kind<Path> ENTRY_MODIFY
10
11     public static final WatchEvent.Kind<Object> OVERFLOW
12 }
```

JDK java.nio.file: доступные события

```
1 package java.nio.file;
2
3 public final class StandardWatchEventKinds {
4
5     public static final WatchEvent.Kind<Path> ENTRY_CREATE
6
7     public static final WatchEvent.Kind<Path> ENTRY_DELETE
8
9     public static final WatchEvent.Kind<Path> ENTRY_MODIFY
10
11     public static final WatchEvent.Kind<Object> OVERFLOW
12 }
```

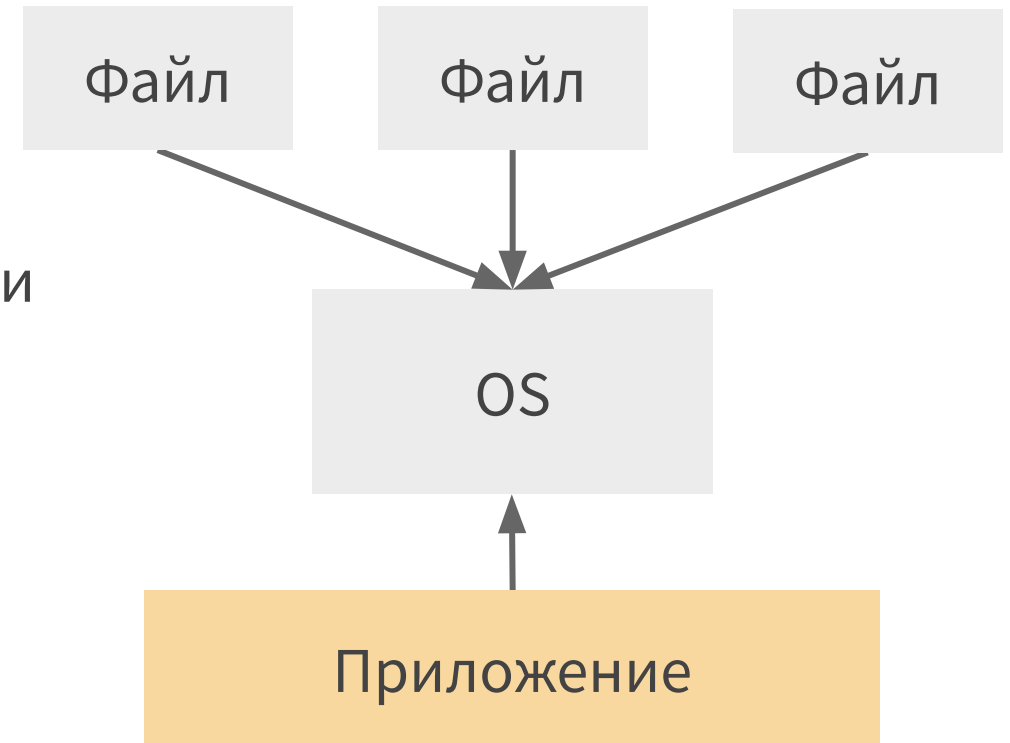
JDK java.nio.file: доступные события

```
1 package java.nio.file;
2
3 public final class StandardWatchEventKinds {
4
5     public static final WatchEvent.Kind<Path> ENTRY_CREATE
6
7     public static final WatchEvent.Kind<Path> ENTRY_DELETE
8
9     public static final WatchEvent.Kind<Path> ENTRY_MODIFY
10
11     public static final WatchEvent.Kind<Object> OVERFLOW
12 }
```

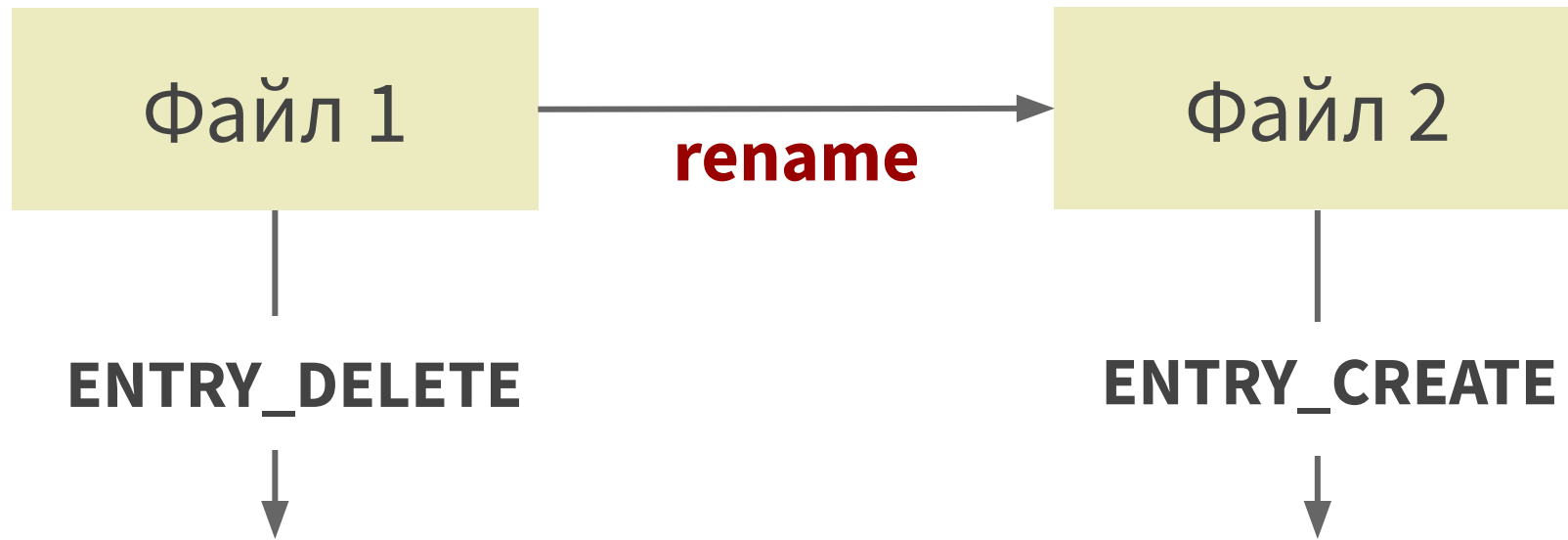
Получение уведомлений O/S: JDK java.nio.file

Недостатки:

- Скучное множество доступных эвентов
- Игнорируется информация о перемещении



JDK java.nio.file: переименование файла



JDK java.nio.file: переименование файла

LinuxWatchService.c *

```
1  JNIEXPORT jintArray JNICALL
2  Java_sun_nio_fs_LinuxWatchService_eventOffsets(...)
3  {
4      jint arr[5];
5      arr[0] = (jint)offsetof(struct inotify_event, wd);
6      arr[1] = (jint)offsetof(struct inotify_event, mask);
7      arr[2] = (jint)offsetof(struct inotify_event, cookie);
8      arr[3] = (jint)offsetof(struct inotify_event, len);
9      arr[4] = (jint)offsetof(struct inotify_event, name);
10     ...
11 }
```

cookie

* <https://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/java.base/linux/native/libnio/fs/LinuxWatchService.c#l54>



JDK java.nio.file: переименование файла

```
1  package sun.nio.fs;
2
3  private static class Poller extends AbstractPoller {
4      public static final int SIZEOF_INOTIFY_EVENT = eventSize();
5      public static final int[] offsets           = eventOffsets();
6      public static final int OFFSETOF_WD        = offsets[0];
7      public static final int OFFSETOF_MASK      = offsets[1];
8      public static final int OFFSETOF_LEN       = offsets[3];
9      public static final int OFFSETOF_NAME      = offsets[4];
10     ...
11 }
```

offsets[2]?

Получение уведомлений O/S: JDK java.nio.file

Подмножество
событий O/S

Недоступны
перемещения



Не подходит при необходимости
использования всех эвентов O/S

Получение уведомлений O/S: inotify/Linux

```
1  #include <sys/inotify.h>
2
3  struct inotify_event {
4      int      wd;
5      uint32_t mask;
6      uint32_t cookie;
7      uint32_t len;
8      char      name[];
9  }
10
11 int inotify_init(void);
12 int inotify_add_watch(int fd, const char *pathname, uint32_t mask);
13 int inotify_rm_watch(int fd, int wd);
```



Получение уведомлений O/S: inotify/Linux

```
1  #include <sys/inotify.h>
2
3  struct inotify_event {
4      int      wd;
5      uint32_t mask;
6      uint32_t cookie;
7      uint32_t len;
8      char      name[];
9  }
10
11 int inotify_init(void);
12 int inotify_add_watch(int fd, const char *pathname, uint32_t mask);
13 int inotify_rm_watch(int fd, int wd);
```



Получение уведомлений O/S: inotify/Linux

```
1  #include <sys/inotify.h>
2
3  struct inotify_event {
4      int      wd;
5      uint32_t mask;
6      uint32_t cookie;
7      uint32_t len;
8      char      name[];
9  }
10
11  int inotify_init(void);
12  int inotify_add_watch(int fd, const char *pathname, uint32_t mask);
13  int inotify_rm_watch(int fd, int wd);
```



inotify/Linux: доступные события

inotify.h *

```
1  #define IN_ACCESS          0x00000001
2  #define IN_MODIFY         0x00000002
3  #define IN_ATTRIB         0x00000004
4  #define IN_CLOSE_WRITE    0x00000008
5  #define IN_CLOSE_NOWRITE  0x00000010
6  #define IN_OPEN           0x00000020
7  #define IN_MOVED_FROM     0x00000040
8  #define IN_MOVED_TO       0x00000080
9  #define IN_CREATE         0x00000100
10 #define IN_DELETE         0x00000200
11 #define IN_DELETE_SELF     0x00000400
12 #define IN_MOVE_SELF      0x00000800
```

* <https://elixir.bootlin.com/linux/v5.3/source/include/uapi/linux/inotify.h#L29>



inotify/Linux: доступные события

inotify.h *

```
1  #define IN_ACCESS          0x00000001
2  #define IN_MODIFY         0x00000002
3  #define IN_ATTRIB         0x00000004
4  #define IN_CLOSE_WRITE    0x00000008
5  #define IN_CLOSE_NOWRITE  0x00000010
6  #define IN_OPEN           0x00000020
7  #define IN_MOVED_FROM     0x00000040
8  #define IN_MOVED_TO       0x00000080
9  #define IN_CREATE         0x00000100
10 #define IN_DELETE         0x00000200
11 #define IN_DELETE_SELF    0x00000400
12 #define IN_MOVE_SELF      0x00000800
```

* <https://elixir.bootlin.com/linux/v5.3/source/include/uapi/linux/inotify.h#L29>



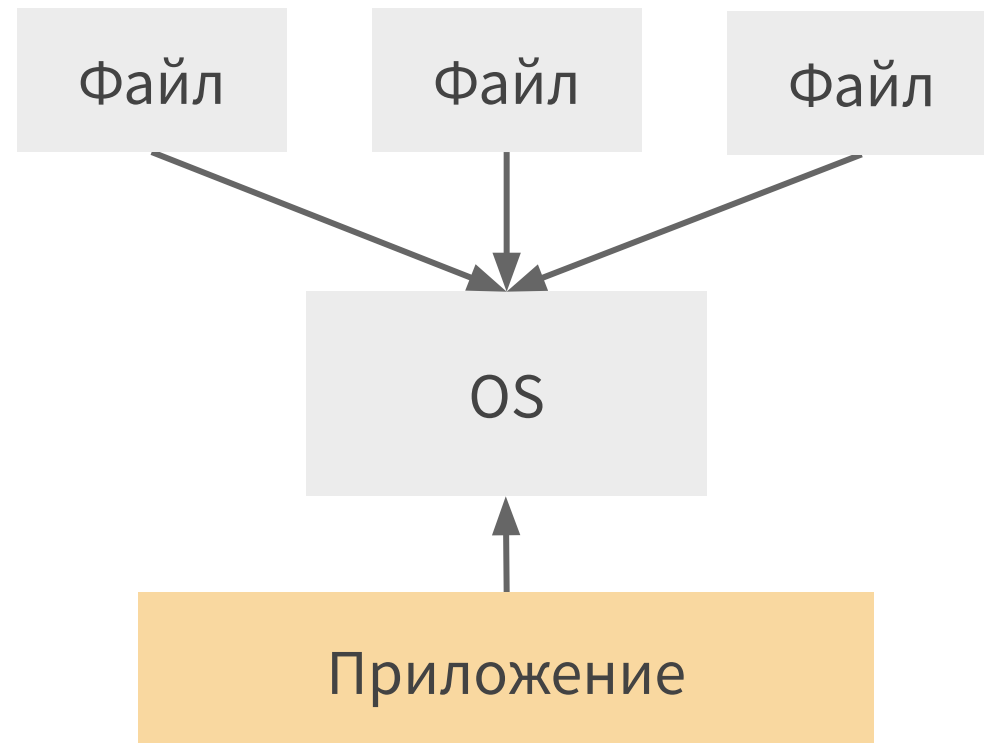
Получение уведомлений O/S: inotify/Linux

Преимущества:

- Полный набор событий OS Linux

Недостатки:

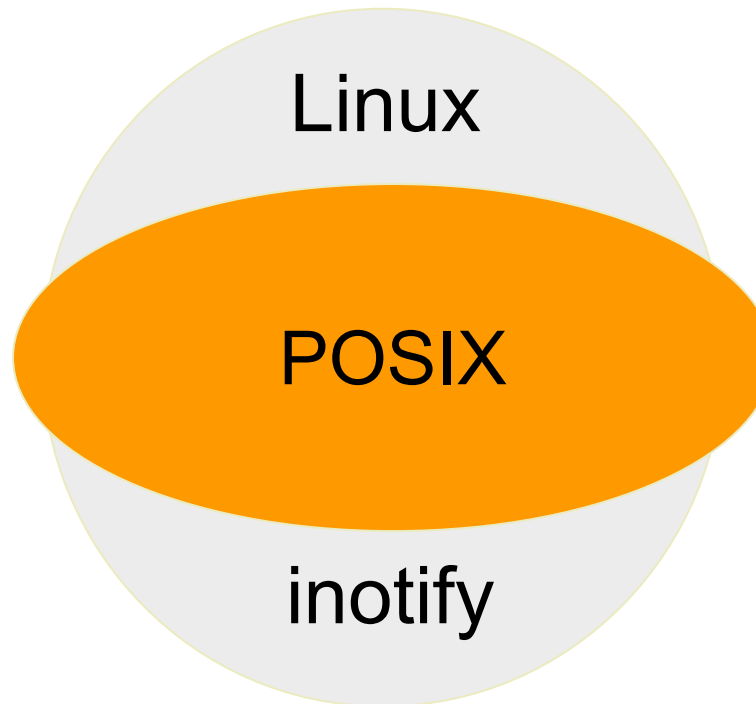
- Не кросс-платформенно / Не POSIX



Стандарт POSIX: описание

Документация POSIX*:

POSIX.1-2017 defines a standard operating system interface and environment, [...] including utility conventions and **C-language header definitions** [...]



* <https://pubs.opengroup.org/onlinepubs/9699919799/>

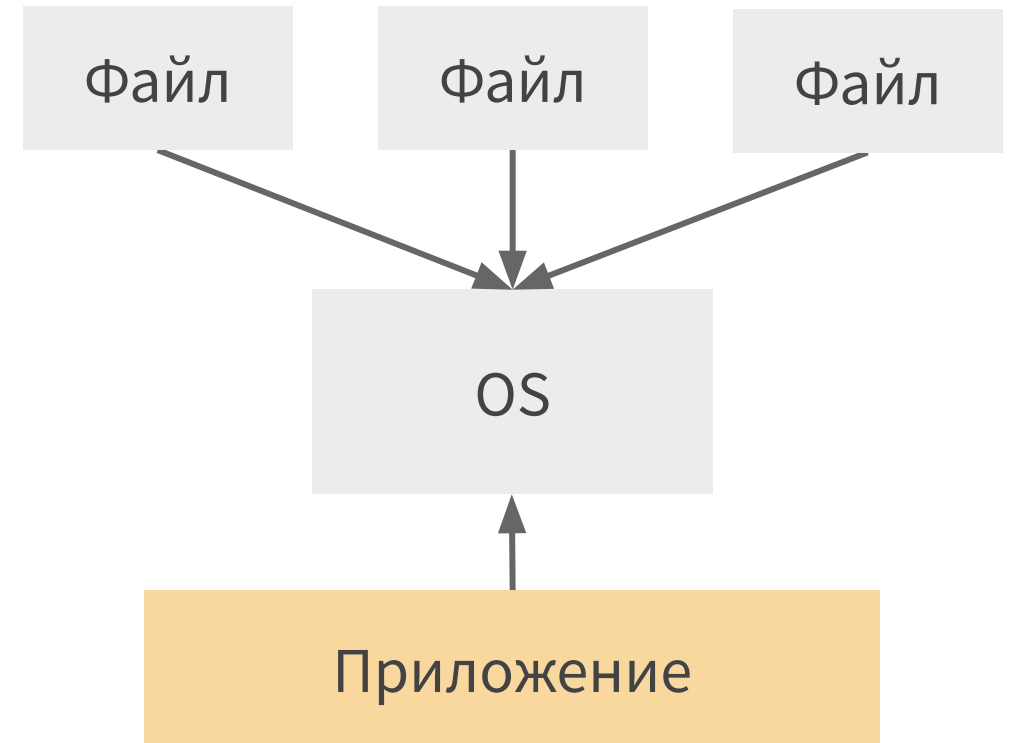
Получение уведомлений O/S: inotify/Linux

Преимущества:

- Полный набор событий OS Linux

Недостатки:

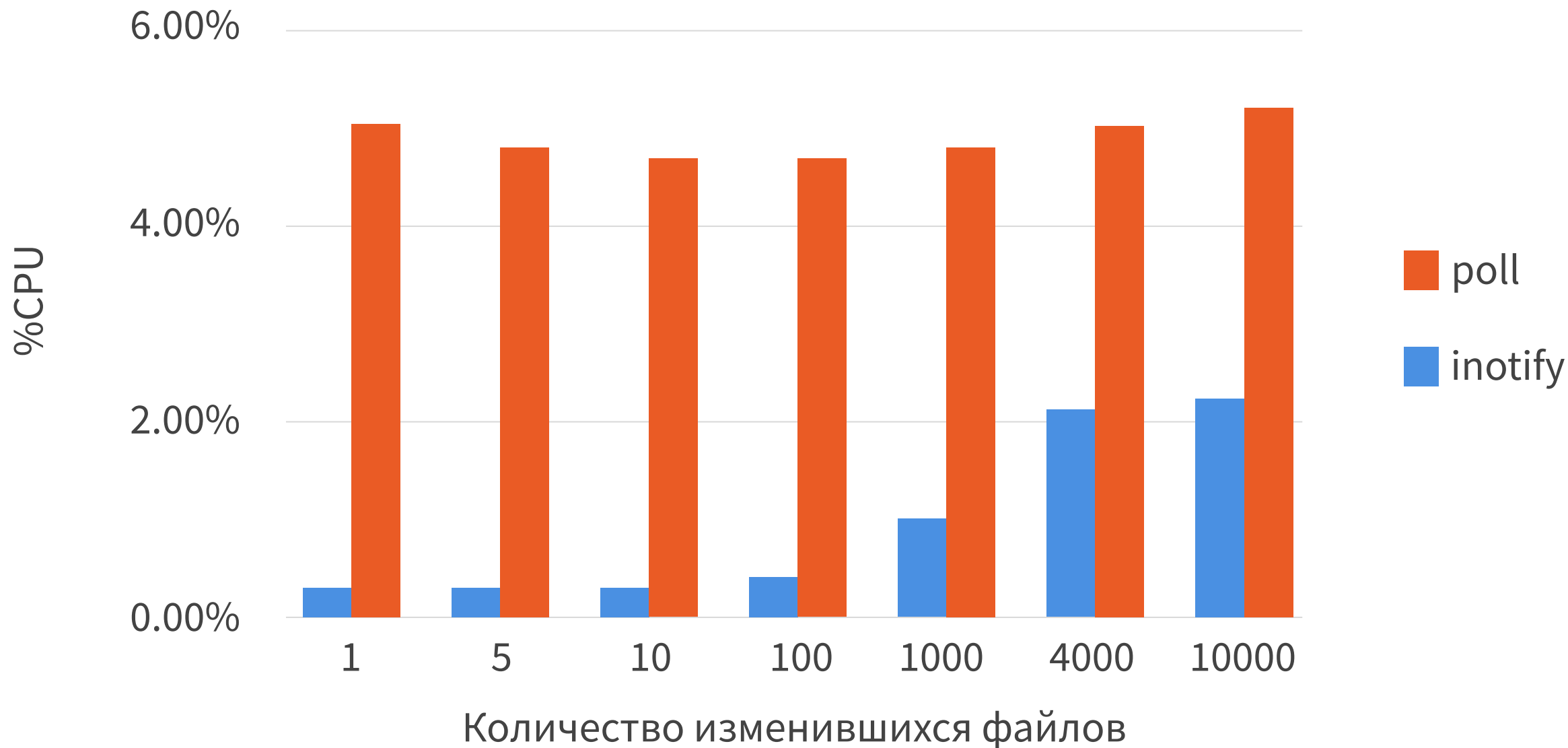
- Не кросс-платформенно / Не POSIX
- Взаимодействие с нативным кодом



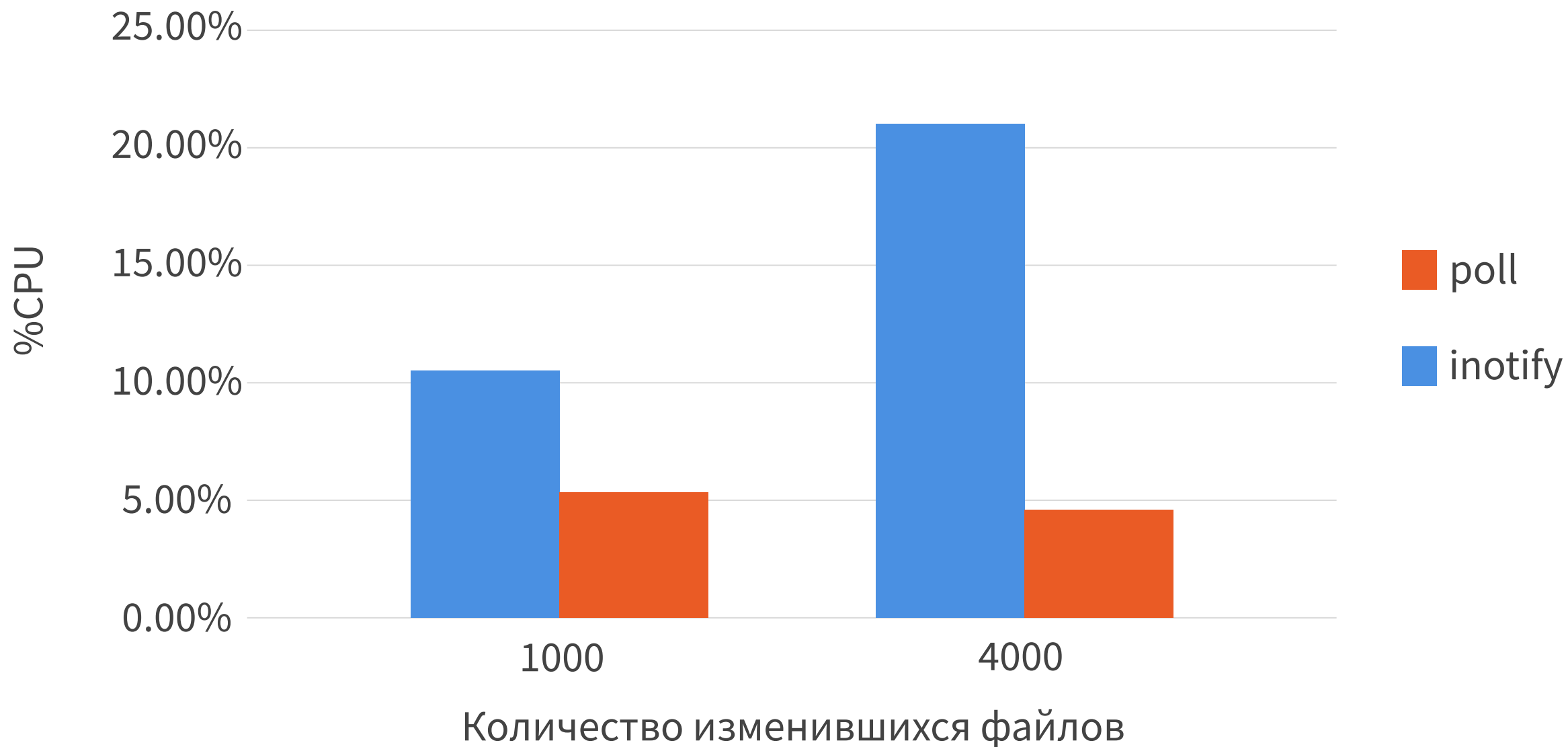
Определение изменившихся файлов: перфоманс

- i7-8550U KbL
- Linux Kernel 5.3.0

Потребление CPU: модификации 1 раз в секунду



Потребление CPU: модификации 10 раз в секунду



Определение изменившихся файлов: резюме

● Поллинг

+ Простота

- Негибкость

- Ресурсоемкость

● JDK WatchService

+ Портруемость

+ Только нужные файлы

- Не все события

● inotify

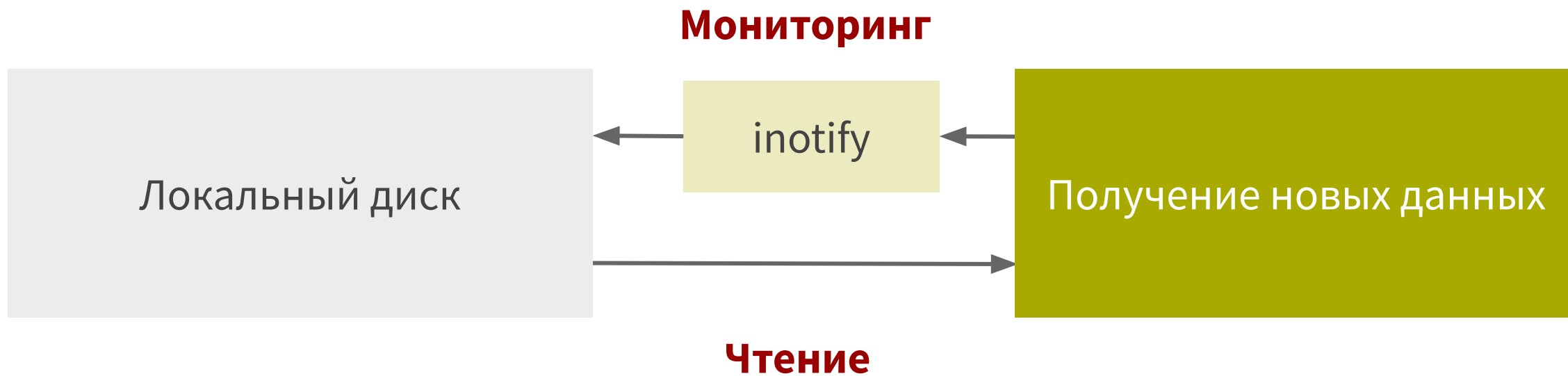
+ Доступны все события

+ Только нужные файлы


- Не POSIX

- Нативный код

Система стриминга: файловое I/O



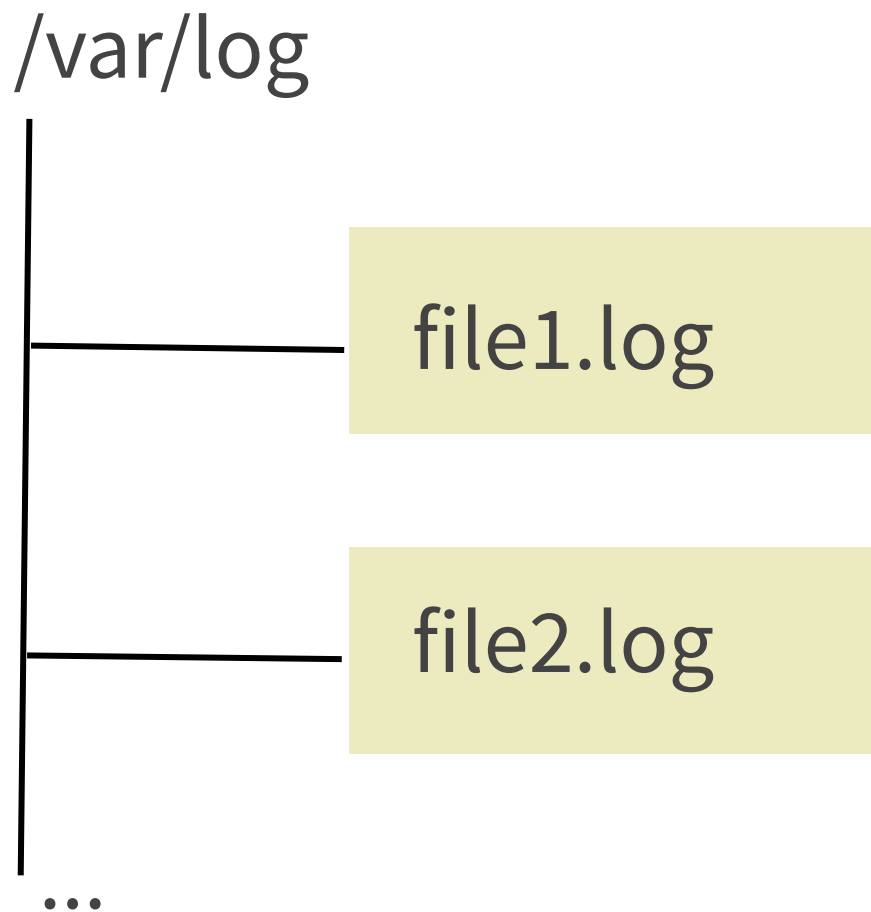
Система стриминга: Файловое I/O

- Определение изменившихся файлов 
- **Определение предыдущей позиции стриминга**
- Чтение новых данных

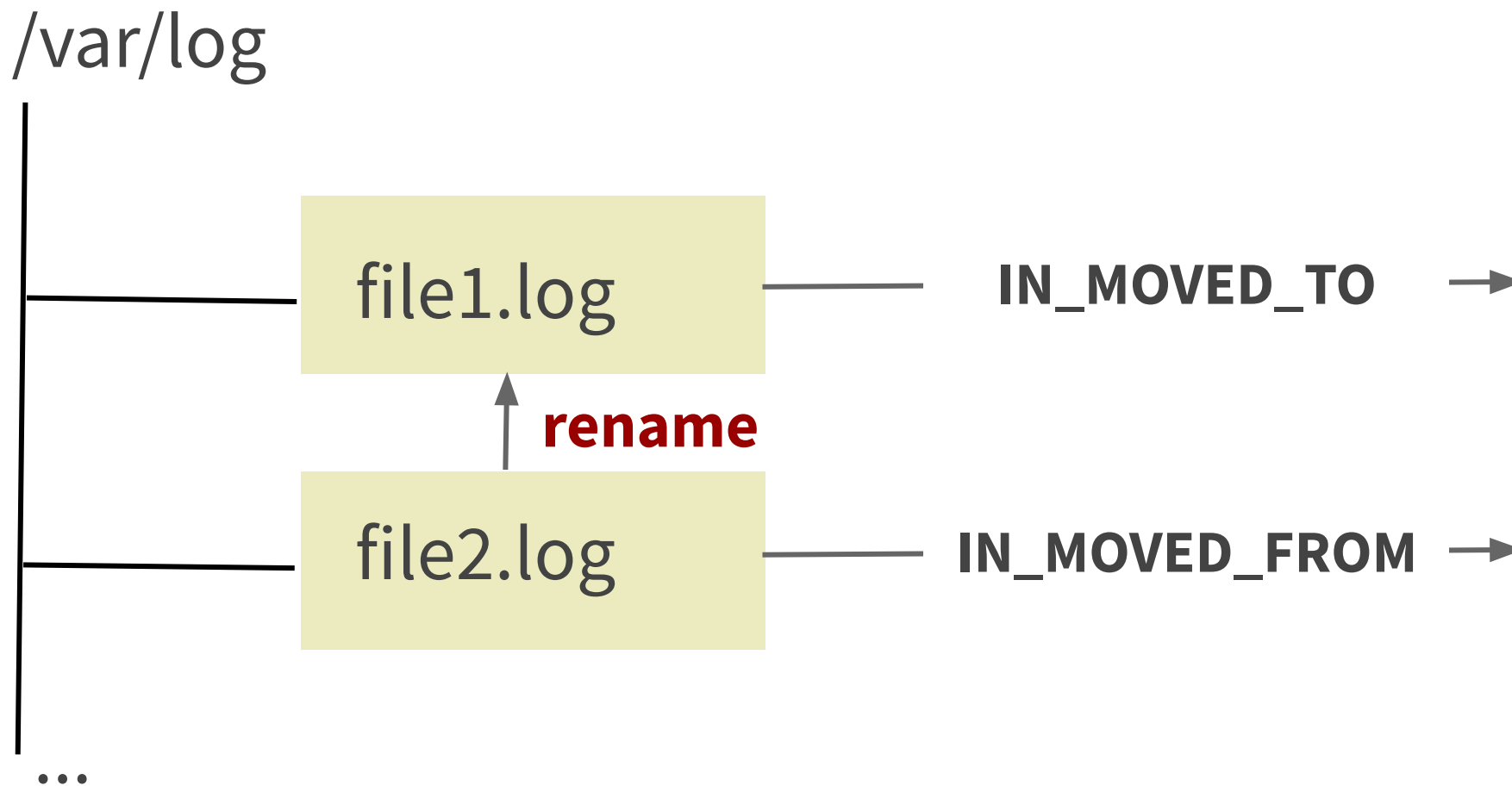
Определение предыдущей позиции стриминга



Идентификация файла: абсолютный путь



Идентификация файла: абсолютный путь

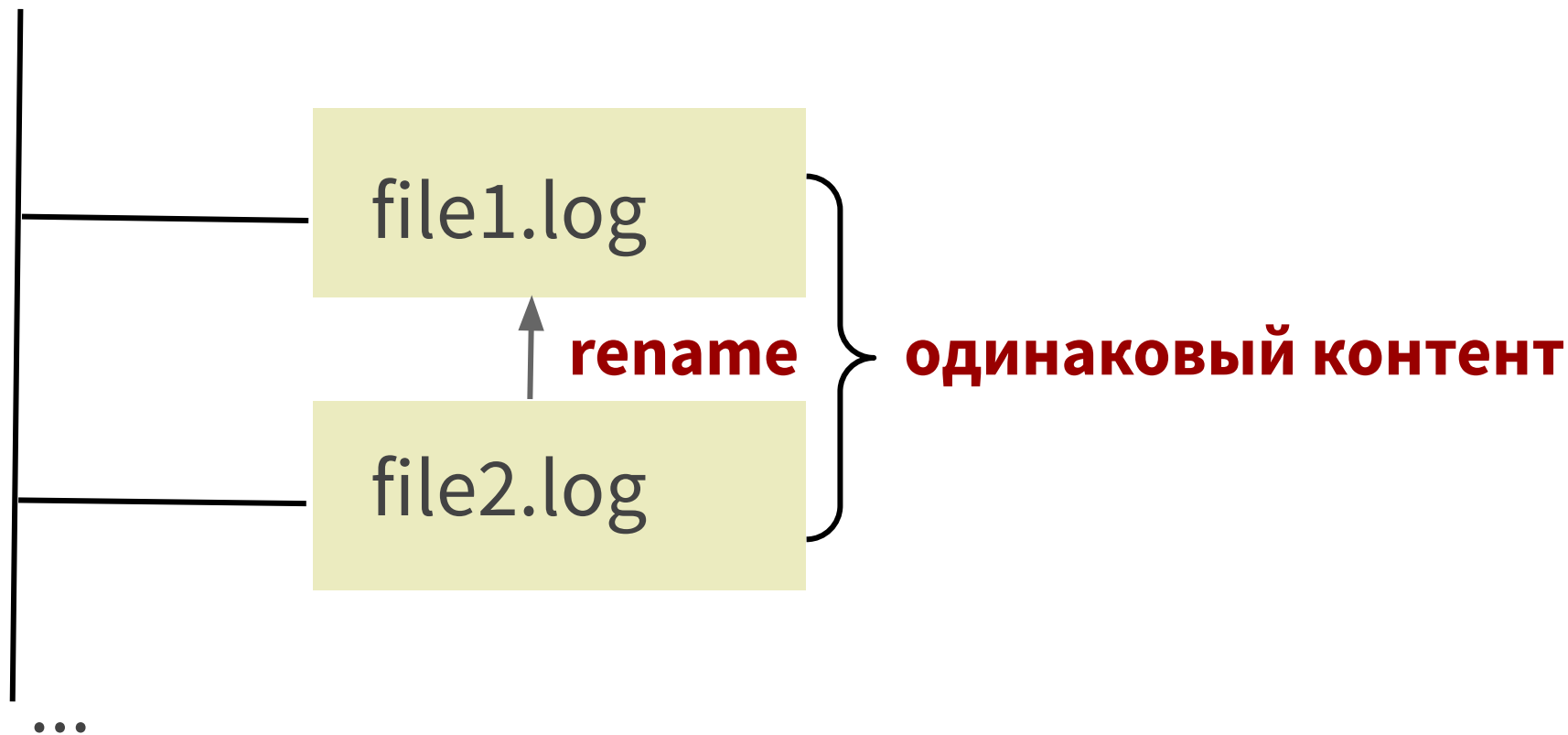


Идентификация файла: хэш первых n байт

 `/var/log/file.log`

Идентификация файла: хэш первых n байт

/var/log



Идентификация файла: POSIX file serial number

Документация POSIX*:

3.176 File Serial Number

A per-file system unique identifier for a file.

* https://pubs.opengroup.org/onlinepubs/9699919799/basedefs/V1_chap03.html

Идентификация файла: POSIX file serial number

Документация POSIX*:

The `<sys/stat.h>` header shall define the structure of the data returned by the functions [`fstat\(\)`](#), [`lstat\(\)`](#), and [`stat\(\)`](#).

```
1  struct stat {  
2      ...  
3      ino_t st_ino;  
4      ...  
5  }
```

* <https://pubs.opengroup.org/onlinepubs/009695399/basedefs/sys/stat.h.html>



POSIX file serial number: JDK java.nio.file

```
1  package java.nio.file;  
2  
3  public final class Files {  
4  
5      public static Object getAttributes(Path path, String attribute,  
6                                          LinkOption... options)  
7  
8  }
```

JDK java.nio.file: пример

```
1 Long inodeNumber = (Long) Files.getAttributes(path, "unix:ino");
```



POSIX file serial number: JDK java.nio.file

UnixNativeDispatcher.c: *

```
1  JNIEXPORT void JNICALL
2  Java_sun_nio_fs_UnixNativeDispatcher_stat0(...) {
3      ...
4      RESTARTABLE(stat64(path, &buf), err);
5      ...
6  }
```

stat64 = stat Ha x86_64

* <https://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/java.base/unix/native/libnio/fs/UnixNativeDispatcher.c#l498>

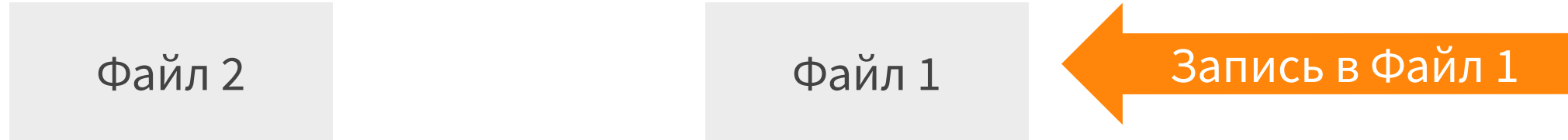


POSIX file serial number: JDK java.nio.file

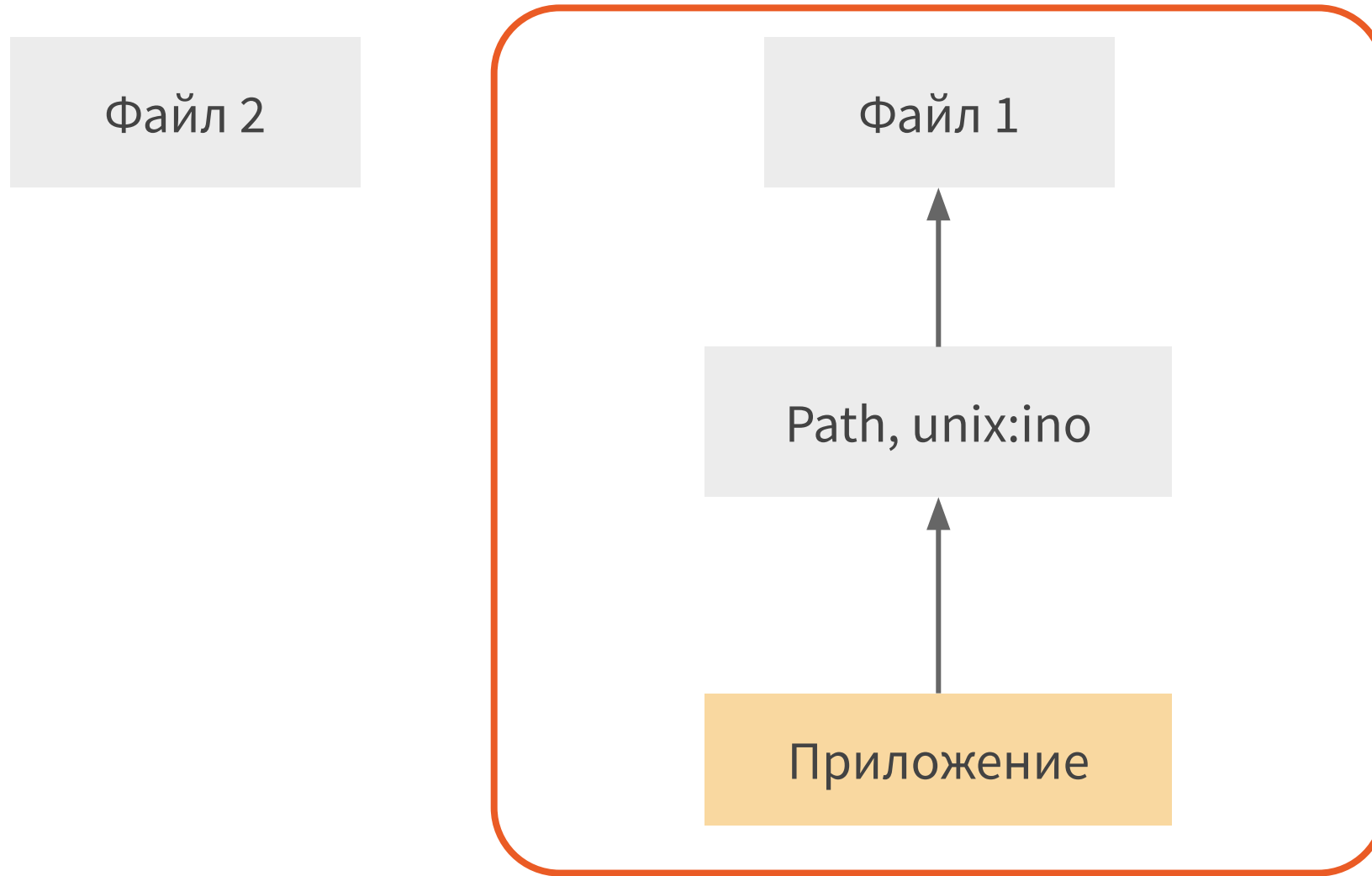
Недостатки:

- Привязка к POSIX
- **Используется путь к файлу**

JDK java.nio.file: Race Condition



JDK java.nio.file: Race Condition

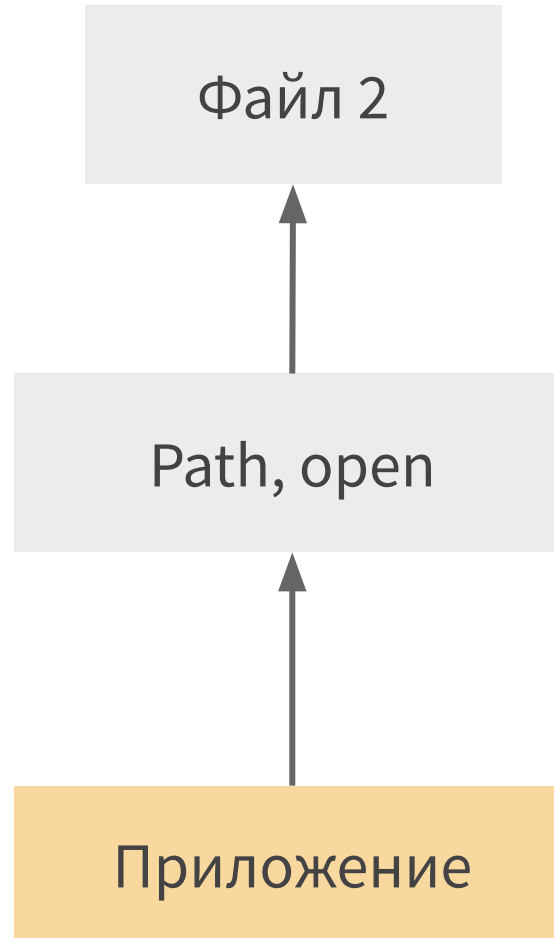


JDK java.nio.file: Race Condition

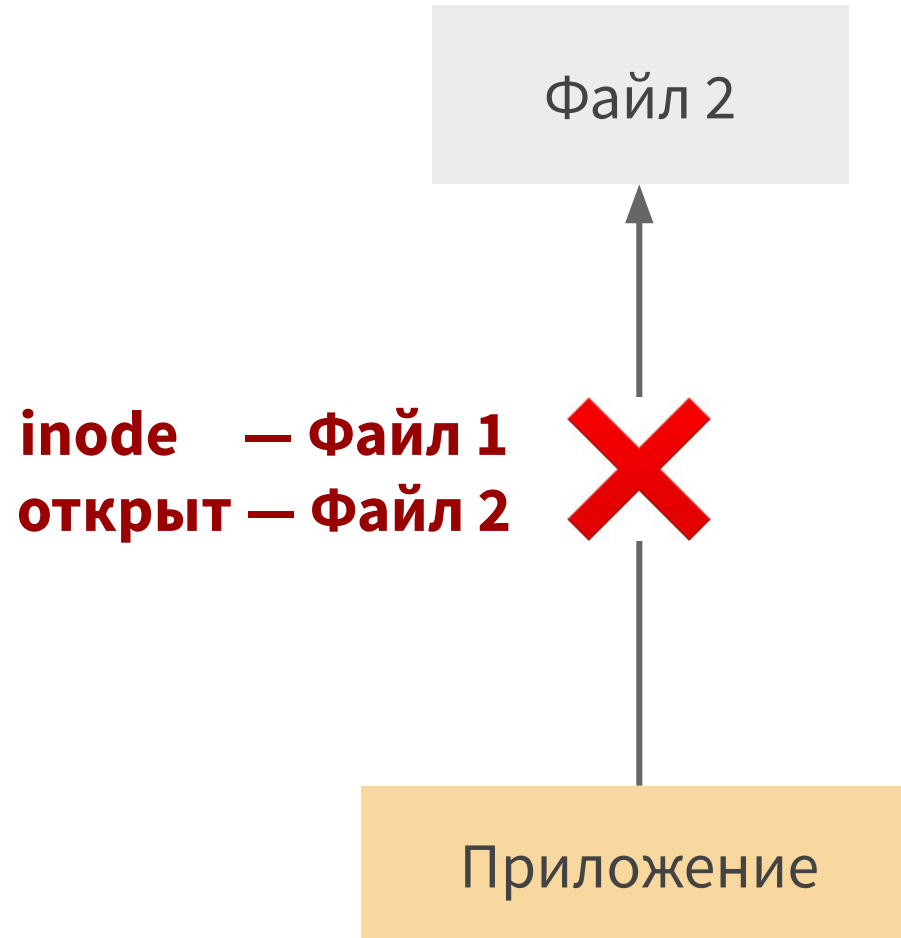


Приложение

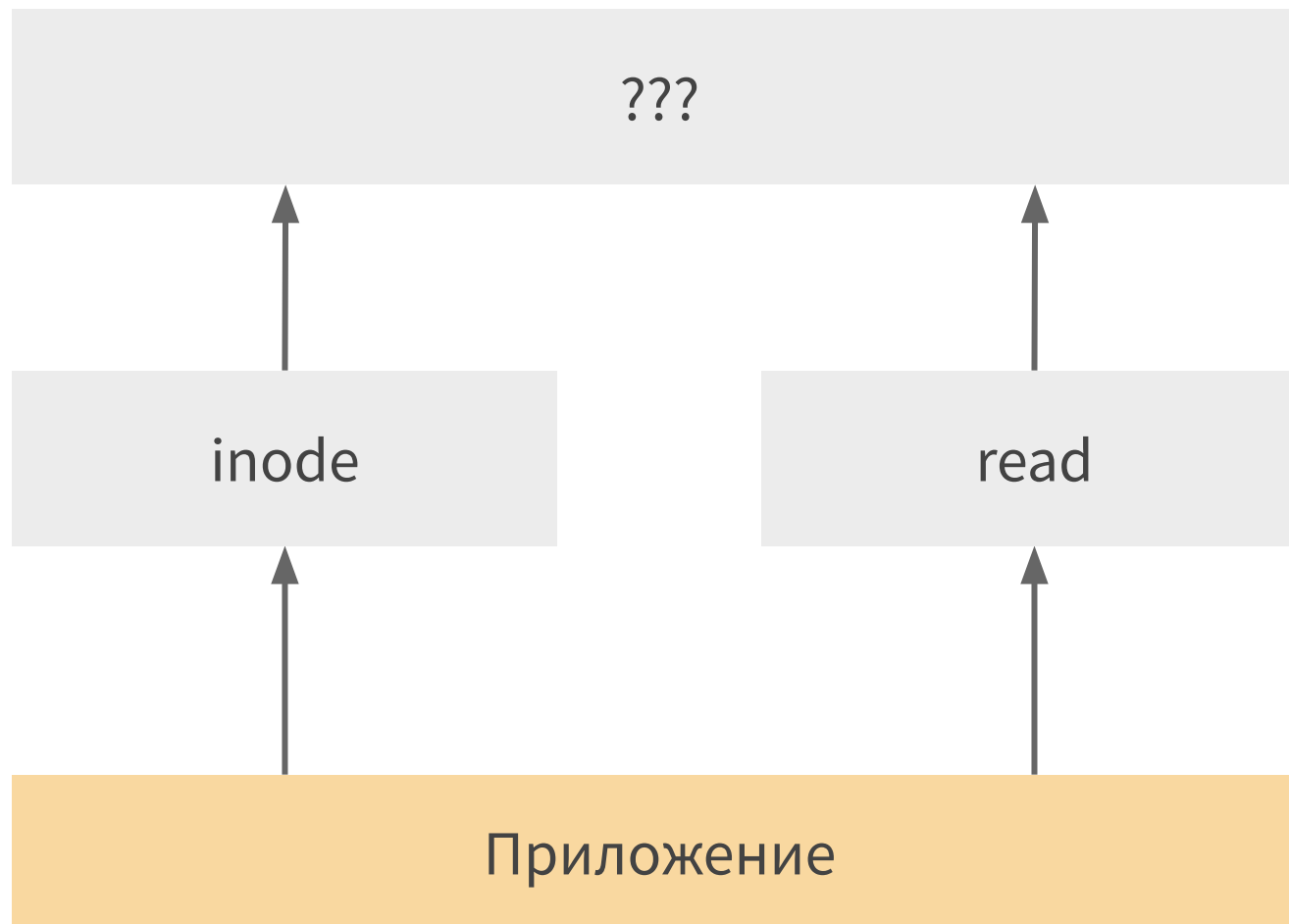
JDK java.nio.file: Race Condition



JDK java.nio.file: Race Condition



Идентификация файла: Race-Free way



POSIX file serial number: POSIX fstat

Документация POSIX*:

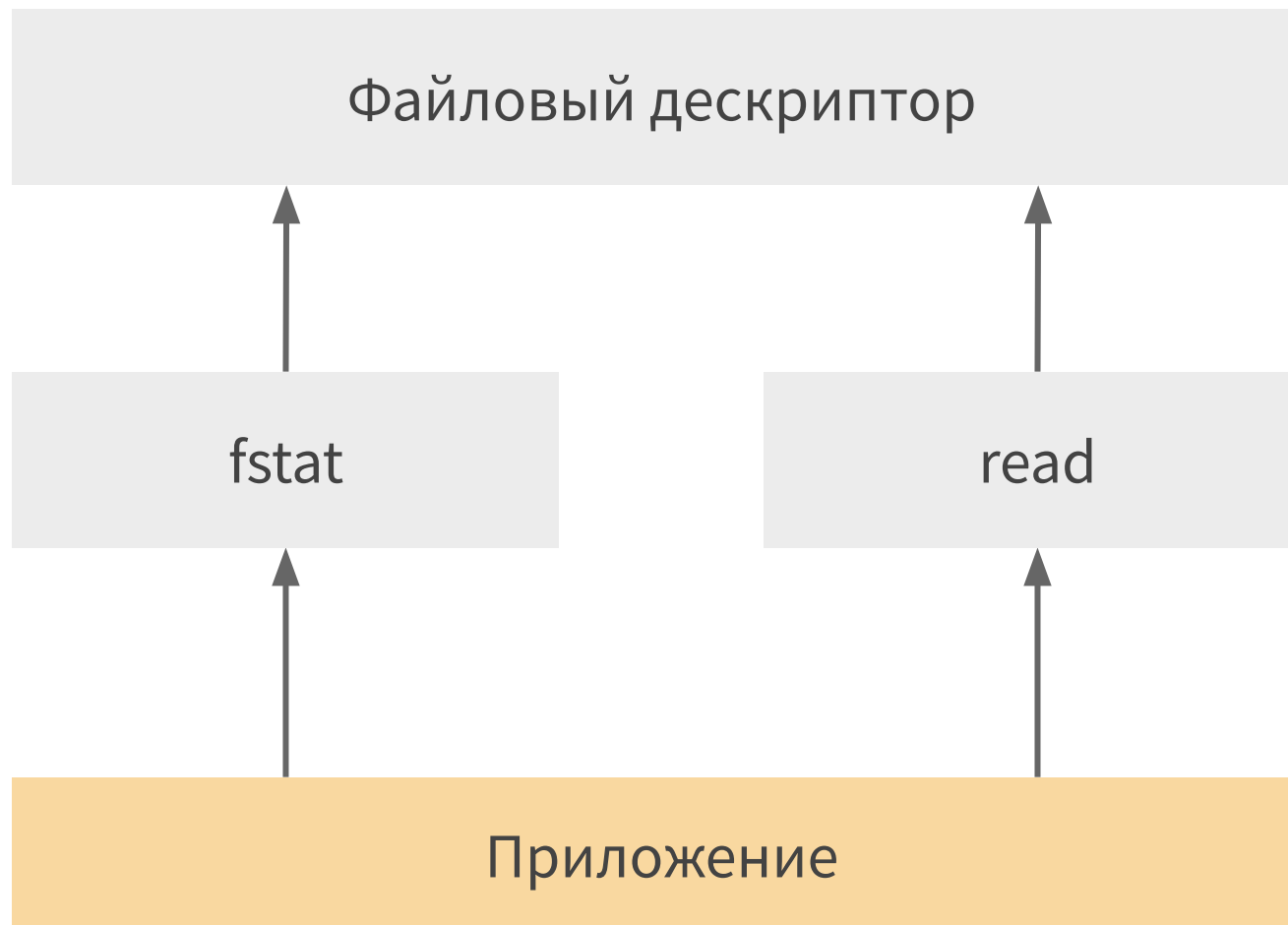
```
1  #include <sys/stat.h>
2
3  int stat(const char *pathname, ...);
4
5  int fstat(int fd, ...);
```

файловый дескриптор

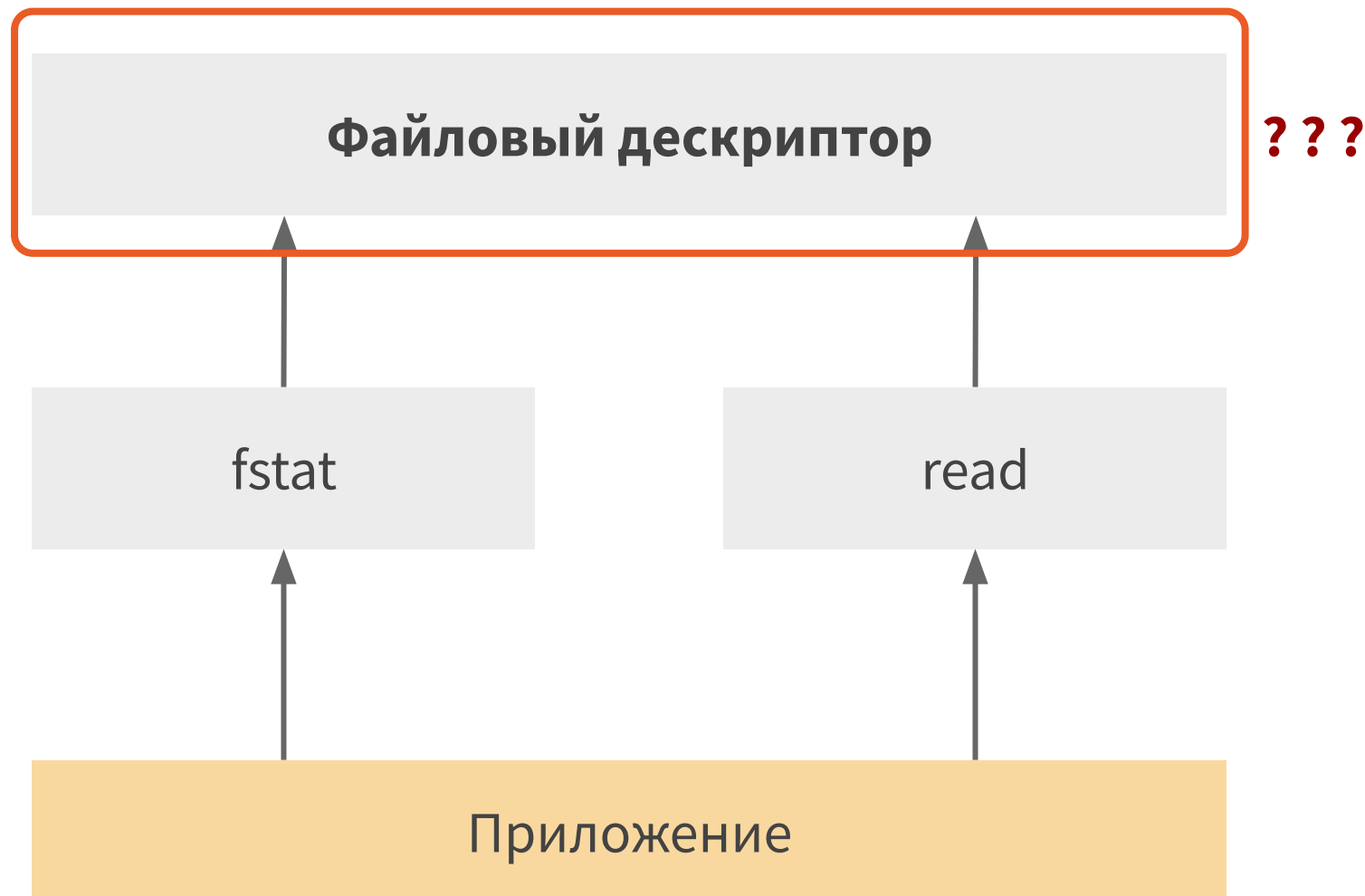
* <https://pubs.opengroup.org/onlinepubs/009695399/functions/fstat.html>



Идентификация файла: Race-Free way



Идентификация файла: Race-Free way



POSIX fstat: файловый дескриптор

Приватное поле:

```
1  package java.io;  
2  
3  public final class FileDescriptor {  
4      private int fd;  
5  }
```

POSIX fstat: файловый дескриптор

Реализация в one-nio:

```
1  package one.nio.os;
2
3  public final class Mem {
4      private static final Field fdField = ...
5
6      public static int getFd(FileDescriptor fd){
7          try {
8              return fdField.getInt(fd);
9          } catch (IllegalAccessException e) {
10             ...
11         }
12     }
13 }
```

POSIX fstat: файловый дескриптор

Нативный метод:

```
1  package com.company;  
2  
3  public class PosixUtils {  
4      public static native int openR0(String path);  
5  }
```

POSIX fstat: файловый дескриптор

Реализация:

```
1  JNIEXPORT void JNICALL
2  Java_com_company_PosixUtils_openR0(...) {
3      const char *native_path = (*env)->GetStringUTFChars(...);
4      int fd = open(native_path, O_RDONLY);
5      (*env)->ReleaseStringUTFChars(...);
6      return fd;
7  }
```



POSIX file serial number: POSIX fstat

Преимущества:

- Отсутствие Race Condition

Недостатки:

- Требуется файловый дескриптор
- Отсутствует в JDK

Идентификация файла: перфоманс

- горячий dentry cache
- горячий page cache

Идентификация файла: перфоманс

```
1  #include <x86intrin.h>
2
3  unsigned long long __rdtsc(void);
4
5  #include <immintrin.h>
6
7  void _mm_clflush(void const *A);
```



Идентификация файла: перфоманс

```
1  #include <x86intrin.h>
2
3  unsigned long long __rdtsc(void);
4
5  #include <immintrin.h>
6
7  void _mm_clflush(void const *A);
```

Чтение Core Time-Stamp Counter



Time-Stamp Counter

Intel System Programming Manual/17.17:

the time-stamp counter increments at a constant rate[...]

The invariant TSC will run at a constant rate in all ACPI P-, C-. and T-states.

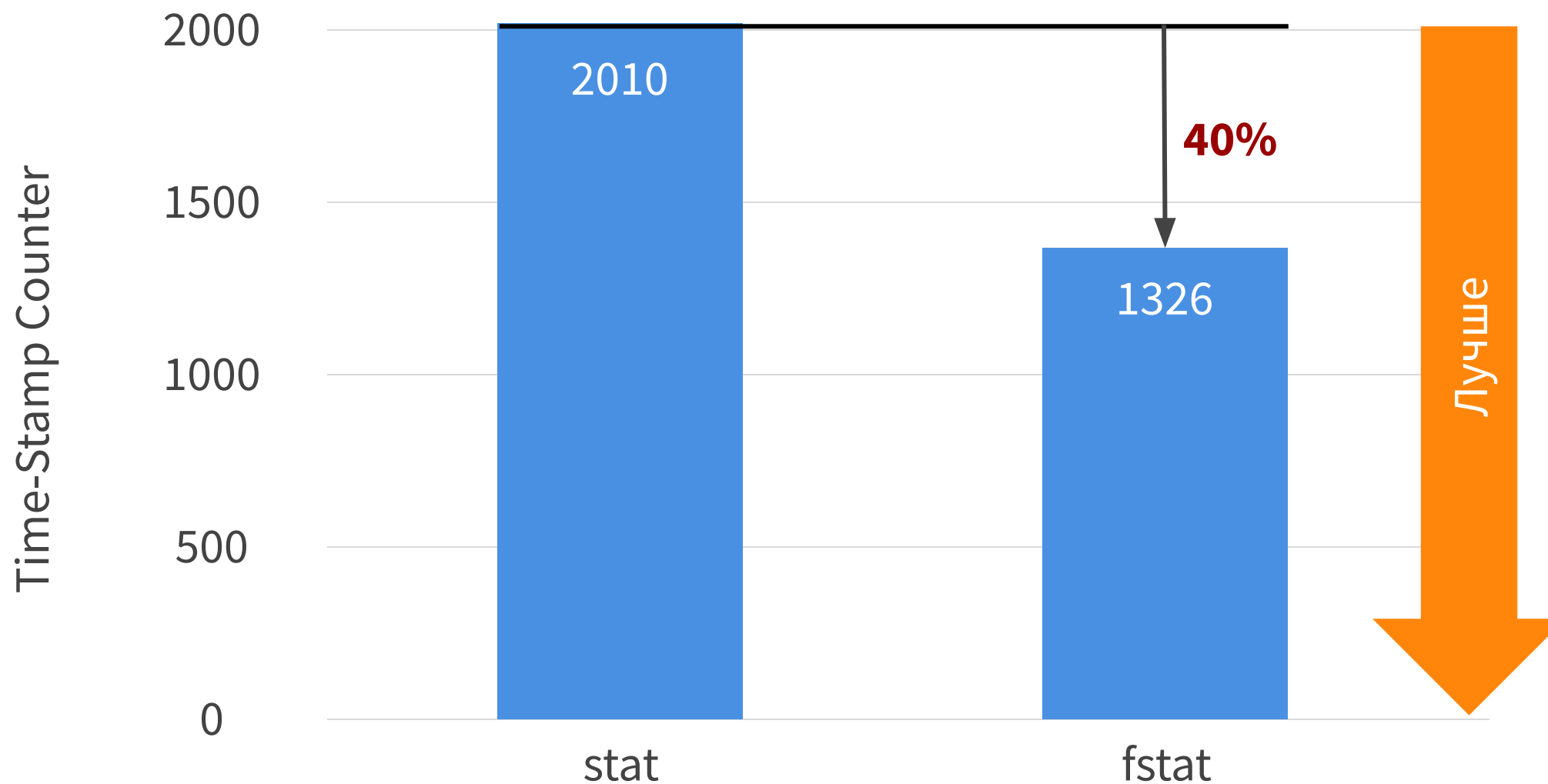
Идентификация файла: перфоманс

```
1  #include <x86intrin.h>
2
3  unsigned long long __rdtsc(void);
4
5  #include <immintrin.h>
6
7  void _mm_clflush(void const *A);
```

Сброс L1/L2/LLC/... по адресу A



Перфоманс: stat vs fstat

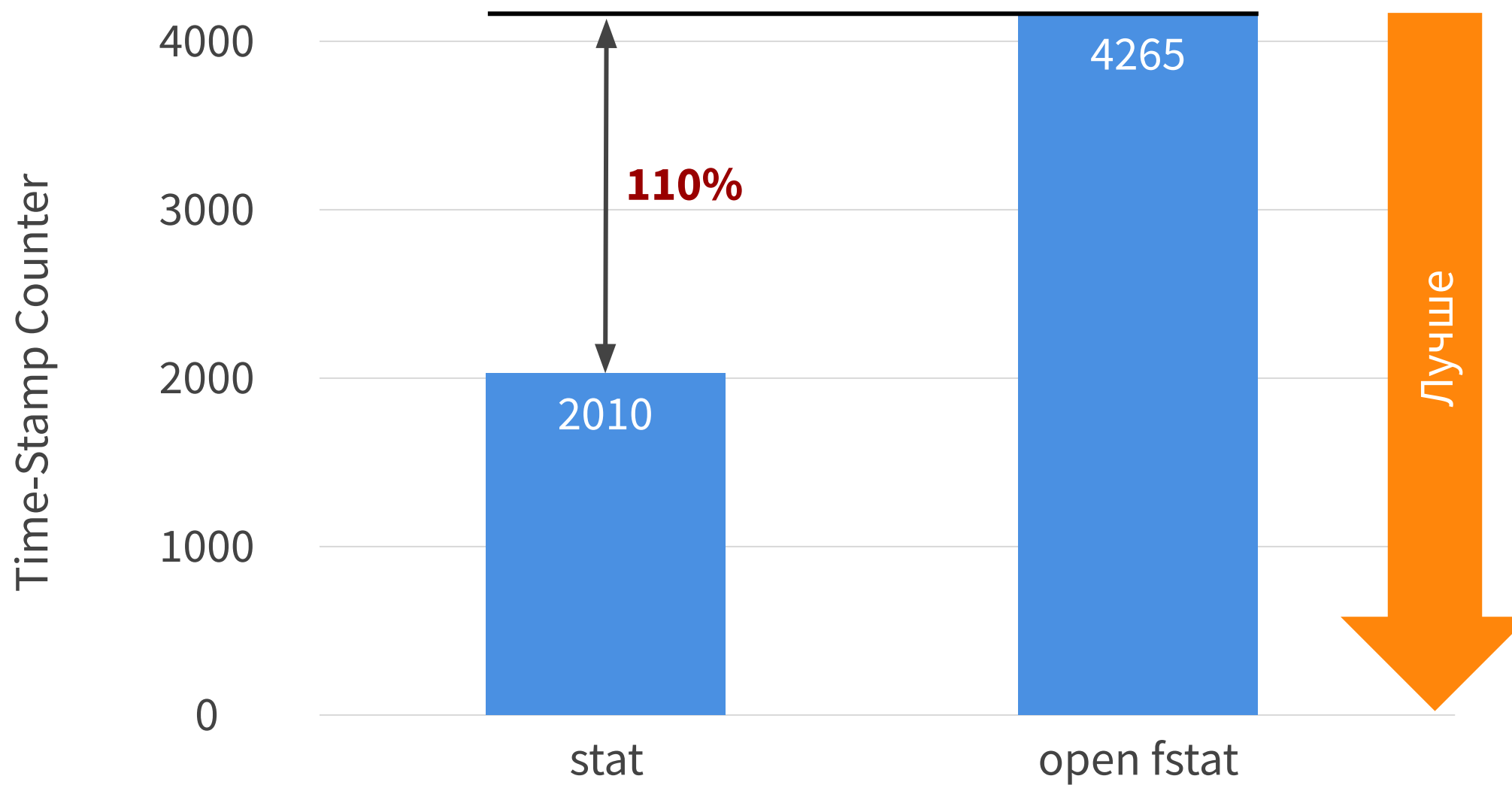


Перфоманс: stat vs fstat

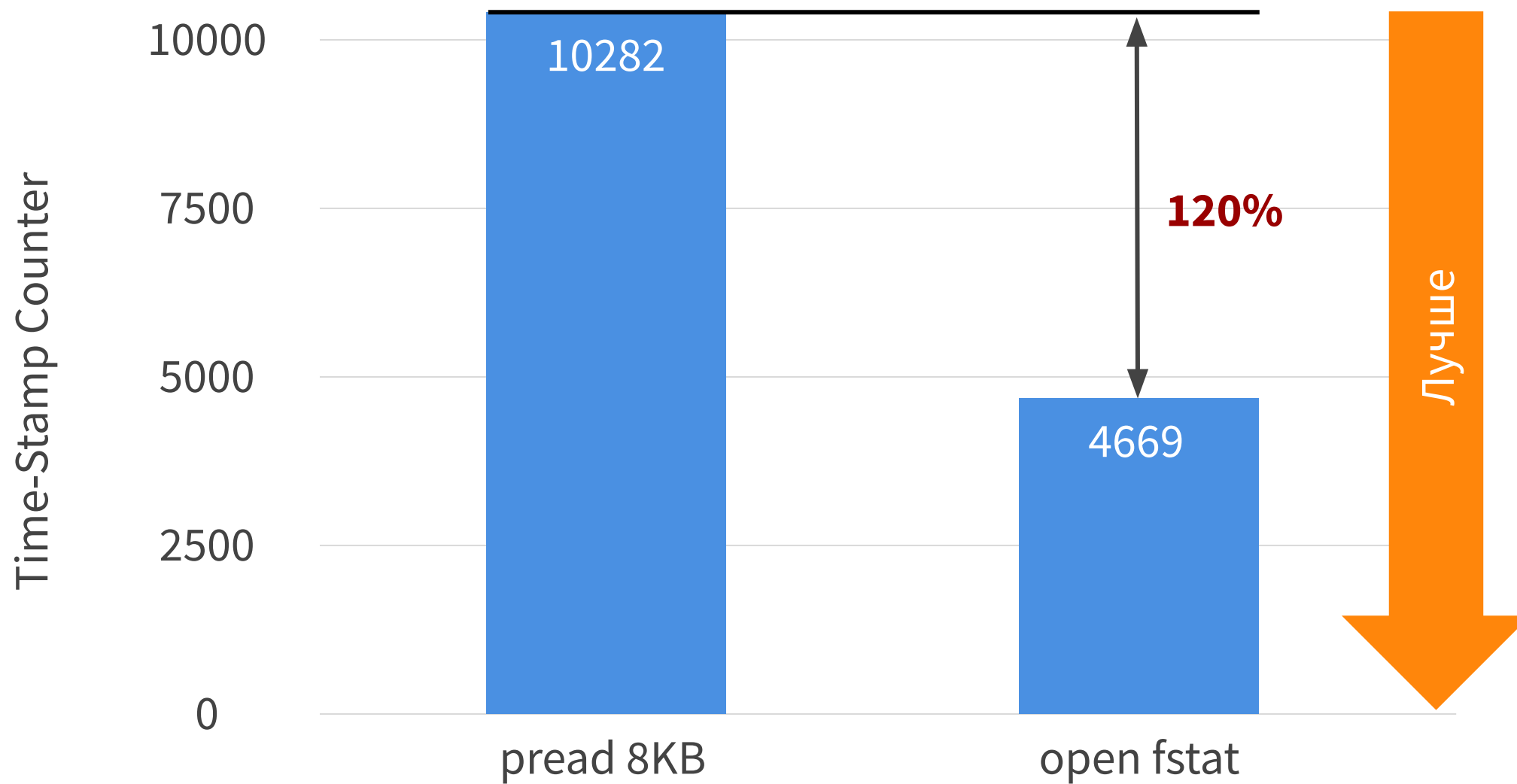
```
- 42,76% 0,28% bin [kernel.kallsyms] [k] user_path_at_empty
- 42,48% user_path_at_empty
  - 30,21% filename_lookup
    - 28,39% path_lookupat
      - 20,55% link_path_walk.part.33
        + 9,64% inode_permission
        + 7,29% walk_component
      + 2,49% complete_walk
      + 2,39% walk_component
      + 1,72% path_init
    + 0,97% putname
  + 11,96% getname_flags
```

Обход каждого из компонентов пути

Перфоманс: stat vs open fstat



Перфоманс: pread vs open fstat, чистый кэш



Идентификация файла: резюме

- **Абсолютный путь**

+ Доступен в inotify

- Не идентифицирует

- **Хэш первых n байт**

+ Портiruемость

- Не идентифицирует

- Перфоманс

- **Inode number**

+ Уникален на FS

+ Перфоманс

- Часть POSIX API

Идентификация файла: резюме

- **JDK Files.getAttribute**

+ Часть JDK

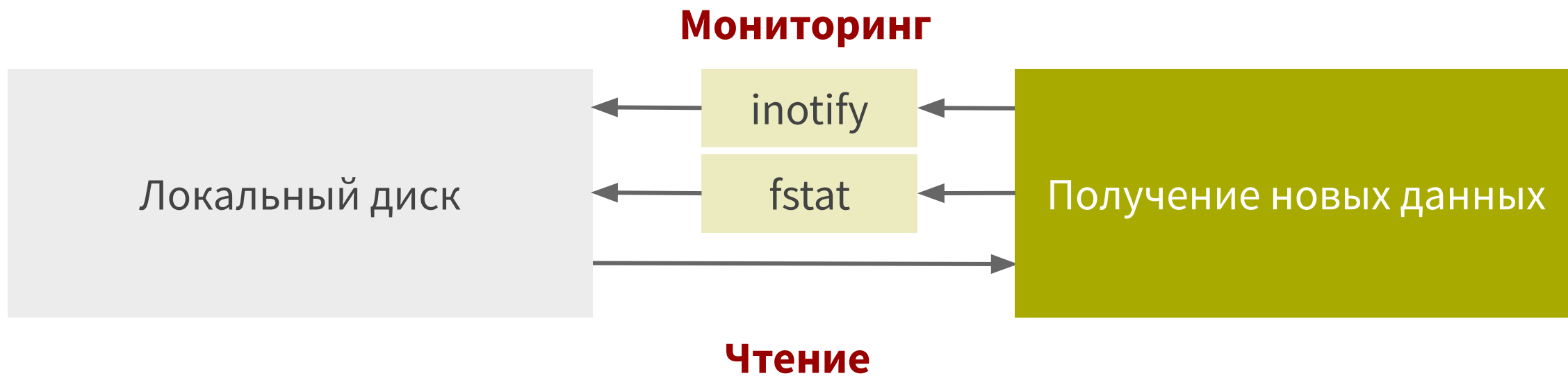
- Race Condition при
использовании пути к файлу

- **POSIX fstat**



+ Не подвержен Race Condition

- Требуется файловый дескриптор
- Недоступен в JDK

Система стриминга: Файловое I/O



Система стриминга: Файловое I/O

- Определение изменившихся файлов 
- Определение предыдущей позиции стриминга 
- **Чтение новых данных**

Чтение новых данных

- Пропускная способность
- Использование on-heap byte[]
 - GZIP до JDK11
 - GCP Connector

Чтение новых данных

Zero-copy file

FileChannel.read

java.io.InputStream

Чтение новых данных

Zero-copy file

Чтение новых данных: Zero-copy file

Схема API:

```
1  package java.nio.channels;
2
3  public abstract class FileChannel {
4
5      public abstract long transferTo(long position, long count,
6                                     WritableByteChannel target)
7
8  }
```


Чтение новых данных: Zero-copy file

FileChannelImpl.c: *

```
1  JNIEXPORT void JNICALL
2  Java_sun_nio_ch_FileChannelImpl_transferTo0(...)
3  {
4      ...
5      #if defined(__linux__)
6          jlong n = sendfile64(dstFD, srcFD, &offset, (size_t)count);
7  }
```



* <http://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/java.base/unix/native/libnio/ch/FileChannelImpl.c#l124>



Чтение новых данных

FileChannel.read

Чтение новых данных: FileChannel.read

Схема API:

```
1  package java.nio.channels;  
2  
3  public abstract class FileChannel {  
4  
5      public abstract int read(ByteBuffer dst, long position)  
6  
7  }
```

Чтение новых данных: FileChannel.read

FileDispatcherImpl.c: *

```
1 JNIEXPORT void JNICALL
2 Java_sun_nio_ch_FileDispatcherImpl_pread0(...)
3 {
4     jint fd = fdval(env, fdo);
5     void *buf = (void *)jlong_to_ptr(address);
6
7     return convertReturnVal(env, pread64(fd, buf, len, offset), ...);
8 }
```

чтение по смещению

* <http://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/java.base/unix/native/libnio/ch/FileDispatcherImpl.c#l88>



Чтение новых данных: FileChannel.read

IOUtil.java: *

```
1  package sun.nio.ch;
2
3  public class IOUtil {
4      static int read(ByteBuffer dst, ...){
5          if(dst instanceof DirectBuffer)
6              return readIntoNativeBuffer(...);
7
8          ByteBuffer bb;
9          bb = Util.getTemporaryDirectBuffer(...);
10         int n = readIntoNativeBuffer(bb, ...);
11         ...
12     }
13 }
```

* <http://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/java.base/share/classes/sun/nio/ch/IOUtil.java#l226>; часть нерелевантного кода опущена

Чтение новых данных: FileChannel.read

IOUtil.java: *

```
1  package sun.nio.ch;
2
3  public class IOUtil {
4      static int read(ByteBuffer dst, ...){
5          if(dst instanceof DirectBuffer)
6              return readIntoNativeBuffer(...);
7
8          ByteBuffer bb;
9          bb = Util.getTemporaryDirectBuffer(...);
10         int n = readIntoNativeBuffer(bb, ...);
11         ...
12     }
13 }
```

Чтение во временный DirectBuffer

* <http://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/java.base/share/classes/sun/nio/ch/IOUtil.java#l226>

Чтение новых данных

java.io.InputStream

Чтение новых данных: java.io.InputStream

Схема API:

```
1  package java.io;  
2  
3  public abstract class InputStream implements Closeable {  
4      public int read(byte[] dst, int off, int len)  
5  }
```


Чтение новых данных: java.io.InputStream

io_util.c: *

```
1  jint readBytes(jint len, ...)
2  {
3      char stackBuf[BUF_SIZE];
4      char *buf = NULL;
5
6      if (len == 0) {
7          return 0;
8      } else if (len > BUF_SIZE) {
9          buf = malloc(len);
10     } else {
11         buf = stackBuf;
12     }
13     ...
14     IO_read(buf, ...)
```

← BUF_SIZE = 8192

* https://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/java.base/share/native/libjava/io_util.c



Чтение новых данных: резюме

- **Zero-copy file**

+ Нет userspace-копии

- Передача без архивации

- **FileChannel.read**

+ Перфоманс в offheap

- Копирование onheap

- **InputStream**

+ Читает в byte[]

- Копирование onheap

- Аллокация в C-heap

Чтение новых данных: резюме

- ???

- + Читает в `byte[]`

- ~~- Копирование `onHeap`~~

- ~~- Аллокация в `C heap`~~

Критическая секция JNI

```
1 void * GetPrimitiveArrayCritical(jarray buf, ...);
```

уведомляем GC: сырой
указатель на регион heap'a

Критическая секция JNI

```
1  void * GetPrimitiveArrayCritical(jarray buf, ...);  
2  
3  void ReleasePrimitiveArrayCritical(jarray buf, void * carray, ...);
```

уведомляем GC: завершение
работы с сырым указателем



Критическая секция JNI: HotSpot VM

jni.cpp: *

```
1  jni_GetPrimitiveArrayCritical(...) {  
2      oop a = lock_gc_or_pin_object(...);  
3      ...  
4  }
```

* <https://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/hotspot/share/prims/jni.cpp#l3174>



Критическая секция JNI: HotSpot VM

jni.cpp: *

```
1  static oop lock_gc_or_pin_object(JavaThread* thread, jobject obj){
2      if(Universe::heap()->supports_object_pinning()){
3          const oop o = JNIHandles::resolve_non_null(obj);
4          return Universe::heap()->pin_object(thread, o);
5      } else {
6          GCLocker::lock_critical(thread);
7          ...
8      }
```

* <https://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/hotspot/share/prims/jni.cpp#l3155>



Критическая секция JNI: HotSpot VM

`supports_object_pinning`

Serial	NO
CMS	NO
G1	NO
Parallel Scavenge	NO
Shenandoah	YES
Z	NO

Критическая секция JNI: нативный Java метод

```
1 package com.company;  
2  
3 public class PosixUtils {  
4     public static native int read(int fd, byte buf[]);  
5 }
```

Файловый дескриптор

Критическая секция JNI: реализация

```
1  JNIEXPORT void JNICALL
2  Java_com_company_PosixUtils_read(JNIEnv *env, jclass jc,
3                                     jint fd, jbyteArray buf){
4      size_t sz = (size_t) (*env)->GetArrayLength(env, buf);
5      void *buf_ptr = (*env)->GetPrimitiveArrayCritical(env, buf, NULL);
6      ssize_t bytesRead = read(fd, buf_ptr, sz);
7      (*env)->ReleasePrimitiveArrayCritical(env, buf, buf_ptr, 0);
8      return (jint) bytesRead;
9  }
```



Критическая секция JNI: реализация

```
1  JNIEXPORT void JNICALL
2  Java_com_company_PosixUtils_read(JNIEnv *env, jclass jc,
3                                     jint fd, jbyteArray buf){
4      size_t sz = (size_t) (*env)->GetArrayLength(env, buf);
5      void *buf_ptr = (*env)->GetPrimitiveArrayCritical(env, buf, NULL);
6      ssize_t bytesRead = read(fd, buf_ptr, sz);
7      (*env)->ReleasePrimitiveArrayCritical(env, buf, buf_ptr, 0);
8      return (jint) bytesRead;
9  }
```



Критическая секция JNI: перфоманс

- горячий page cache

Критическая секция JNI: перфоманс

- горячий page cache

	InputStream	criticalRead
2 GB	0.36 sec	0.27 sec
4 GB	0.59 sec	0.47 sec
6 GB	0.95 sec	0.72 sec



20-25 %




Система стриминга: Файловое I/O

Мониторинг



Чтение

Система стриминга: файловое I/O

- Определение изменившихся файлов 
- Определение предыдущей позиции стриминга 
- Чтение новых данных 

Критическая секция JNI: анализ

```
$ sudo perf record --call-graph dwarf -F 9123 -p <app_pid>
```


Критическая секция JNI: анализ


```
$ sudo perf record --call-graph dwarf -F 9123 -p <app_pid>
```

```
- 68,44% 68,36% java [kernel.kallsyms] [k] copy_user_enhanced_fast_string
  68,36% 0x7f50ec3e4519
    Java_com_test_PosixUtils_read
    read (inlined)
    __GI___libc_read (inlined)
    entry_SYSCALL_64_after_hwframe
    do_syscall_64
    __x64_sys_read
    ksys_read
    vfs_read
    __vfs_read
    new_sync_read
    ext4_file_read_iter
    generic_file_read_iter
  - copy_page_to_iter
    68,02% copy_user_enhanced_fast_string
```

Критическая секция JNI: анализ

```
$ sudo perf record --call-graph dwarf -F 9123 -p <app_pid>
```

```
- 68,44% 68,36% java [kernel.kallsyms] [k] copy_user_enhanced_fast_string
  68,36% 0x7f50ec3e4519
    Java_com_test_PosixUtils_read
    read (inlined)
    __GI___libc_read (inlined)
    entry_SYSCALL_64_after_hwframe
    do_syscall_64
    __x64_sys_read
    ksys_read
    vfs_read
    __vfs_read
    new_sync_read
    ext4_file_read_iter
    generic_file_read_iter
  - copy_page_to_iter
    68,02% copy_user_enhanced_fast_string
```

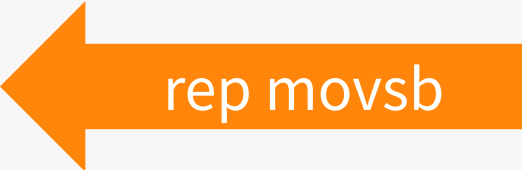


Копирование из ядра пользователю

Копирование памяти в ядре

copy_user_64.S: *

```
1  ENTRY(copy_user_enhanced_fast_string)
2      ...
3      movl %edx,%ecx
4  1:  rep
5      movsb
6      xorl %eax,%eax
7      ASM_CLAC
8      ret
9      ...
```



* https://elixir.bootlin.com/linux/v5.3/source/arch/x86/lib/copy_user_64.S#L200

Копирование памяти в ядре

- Почему не используется AVX?
- Будет ли AVX эффективнее `rep movsb`?
- Если AVX быстрее, то можно ли его использовать?

Копирование памяти в ядре

- Почему не используется AVX?
- Будет ли AVX эффективнее `rep movsb`?
- Если AVX быстрее, то можно ли его использовать?

Копирование памяти в ядре: SysV ABI

Спецификация AMD64 SysV ABI: *

A system-call is done via the syscall instruction. The kernel destroys registers **%rcx** and **%r11**.

* <https://github.com/hjl-tools/x86-psABI/wiki/X86-psABI>

System Call: сохранение и восстановление состояния

calling.h: *

```
1  .macro PUSH_AND_CLEAR_REGS ...  
2      pushq %rsi
```

```
1  .macro POP_REGS ...  
2      popq %rsi
```

* <https://elixir.bootlin.com/linux/latest/source/arch/x86/entry/calling.h>

System Call: сохранение и восстановление состояния

calling.h: *

```
1  .macro PUSH_AND_CLEAR_REGS ...
2      pushq %rsi
3      pushq %rdi
```

```
1  .macro POP_REGS ...
2      popq %rdi
3      popq %rsi
```

* <https://elixir.bootlin.com/linux/latest/source/arch/x86/entry/calling.h>

System Call: сохранение и восстановление состояния

calling.h: *

```
1  .macro PUSH_AND_CLEAR_REGS ...
2      pushq %rsi
3      pushq %rdi
4      pushq %rdx
```

```
1  .macro POP_REGS ...
2      popq %rdx
3      popq %rdi
4      popq %rsi
```

* <https://elixir.bootlin.com/linux/latest/source/arch/x86/entry/calling.h>; аргумент \rdx развернут в %rdx для удобства чтения

System Call: сохранение и восстановление состояния

calling.h: *

```
1  .macro PUSH_AND_CLEAR_REGS ...
2      pushq %rsi
3      pushq %rdi
4      pushq %rdx
5      ...
6      pushq %r13
7      pushq %r14
8      pushq %r15
```

```
1  .macro POP_REGS ...
2      popq %r15
3      popq %r14
4      popq %r13
5      ...
6      popq %rdx
7      popq %rdi
8      popq %rsi
```

* <https://elixir.bootlin.com/linux/latest/source/arch/x86/entry/calling.h>; аргумент \rdx развернут в %rdx для удобства чтения

Копирование памяти в ядре

- Почему не используется AVX?
- Будет ли AVX эффективнее `rep movsb`?
- Если AVX быстрее, то можно ли его использовать?

Сравнение AVX2 и rep movsb

System.arraycopy *:

```
1  else if (UseAVX == 2) {  
2      __ vmovdqu(xmm0, Address(..., -56));  
3      __ vmovdqu(Address(..., -56), xmm0);  
4      __ vmovdqu(xmm1, Address(..., -24));  
5      __ vmovdqu(Address(..., -24), xmm1);  
6  }
```

* https://hg.openjdk.java.net/jdk/jdk12/file/06222165c35f/src/hotspot/cpu/x86/stubGenerator_x86_64.cpp#l1263

Сравнение AVX2 и rep movsb

libc-2.27/memcpy *:

```
1  (gdb) disas __memmove_avx_unaligned_erms
2  ...
3  <+76>:  mov      %rdx,%rcx
4  <+79>:  rep movsb %ds:(%rsi),%es:(%rdi)
5  <+81>:  retq
6  ...
7  <+178>: cmp      $0x1000,%rdx
8  <+185>: ja      0x18eafb
```



rep movsb



4КБ и более

Сравнение memcpy и System.arraycopy: Java memcpy

```
1  public class Memcpy {  
2      public static native void arrayMemcpy(  
3          byte src[], int srcOff,  
4          byte dst[], int dstOff,  
5          int len  
6      );  
7  }
```

Сравнение memcpy и System.arraycopy: Java memcpy

```
1 JNIEXPORT void JNICALL Java_com_company_Memcpy_arrayMemcpy(...) {  
2     ...  
3 }  
4  
5 JNIEXPORT void JNICALL JavaCritical_com_company_Memcpy_arrayMemcpy(...) {  
6     memcpy((char *) dst + dst_off,  
7           (char *) src + src_off,  
8           (size_t) len);  
9 }
```

Сравнение memcpy и System.arraycopy

```
1  @Param({"8192", "131072", "8388608", "16777216", "33554432"})
2  public int size;
3  public byte src[], dst[];
4
5  @Benchmark
6  public void arraycopy(){
7      System.arraycopy(src, 0, dst, 0, size);
8  }
9
10 @Benchmark
11 public void memcpy(){
12     Memcpy.arrayMemcpy(src, 0, dst, 0, size);
13 }
```


Сравнение memcpy и System.arraycopy

	memcpy, μ s	System.arraycopy, μ s
8 KB	0.172	0.132
128 KB	3.439	4.171
8 MB	542.921	959.101
16 MB	1393.213	2209.657
32 MB	3128.653	4702.299

Сравнение memcpy и System.arraycopy

	memcpy, μ s	System.arraycopy, μ s
8 KB	0.172	0.132
128 KB	3.439	4.171
8 MB	542.921	959.101
16 MB	1393.213	2209.657
32 MB	3128.653	4702.299



60-90 %

Non-Temporal Stores

```
1  (gdb) disas __memmove_avx_unaligned_erms
2  <+43>:      cmp      0x261b76(%rip),%rdx
3  <+50>:      jae      0x18ec2d
4  ...
5  <+770>:     vmovntdq %ymm0, (%rdi)
6  <+774>:     vmovntdq %ymm1, 0x20(%rdi)
7  <+779>:     vmovntdq %ymm2, 0x40(%rdi)
8  <+784>:     vmovntdq %ymm3, 0x60(%rdi)
9  ...
10 <+805>:     sfence
```

x86 shared non-temporal threshold

Non-Temporal Stores: x86 non-temporal threshold

```
1  (gdb) p __x86_shared_non_temporal_threshold
2  $1 = 6291356 #6MB
```

Non-Temporal Stores

```
1  (gdb) disas __memmove_avx_unaligned_erms
2  <+43>:      cmp      0x261b76(%rip),%rdx
3  <+50>:      jae      0x18ec2d
4  ...
5  <+770>:     vmovntdq %ymm0, (%rdi)
6  <+774>:     vmovntdq %ymm1, 0x20(%rdi)
7  <+779>:     vmovntdq %ymm2, 0x40(%rdi)
8  <+784>:     vmovntdq %ymm3, 0x60(%rdi)
9  ...
10 <+805>:     sfence
```

vmovntdq - Non-Temporal store

Non-Temporal Stores: свойства

Intel Software Optimization Manual/8.4.1:

- Write combining — Successive writes to the same cache line are combined
- Write collapsing — Successive writes to the same bytes(s) result in only the last write being visible
- Weakly ordered — No ordering is preserved between WC stores or between WC stores and other loads or stores
- **Uncacheable and not write-allocating** — Stored data is written around the cache and will not generate a read-for-ownership bus request for the corresponding cache line

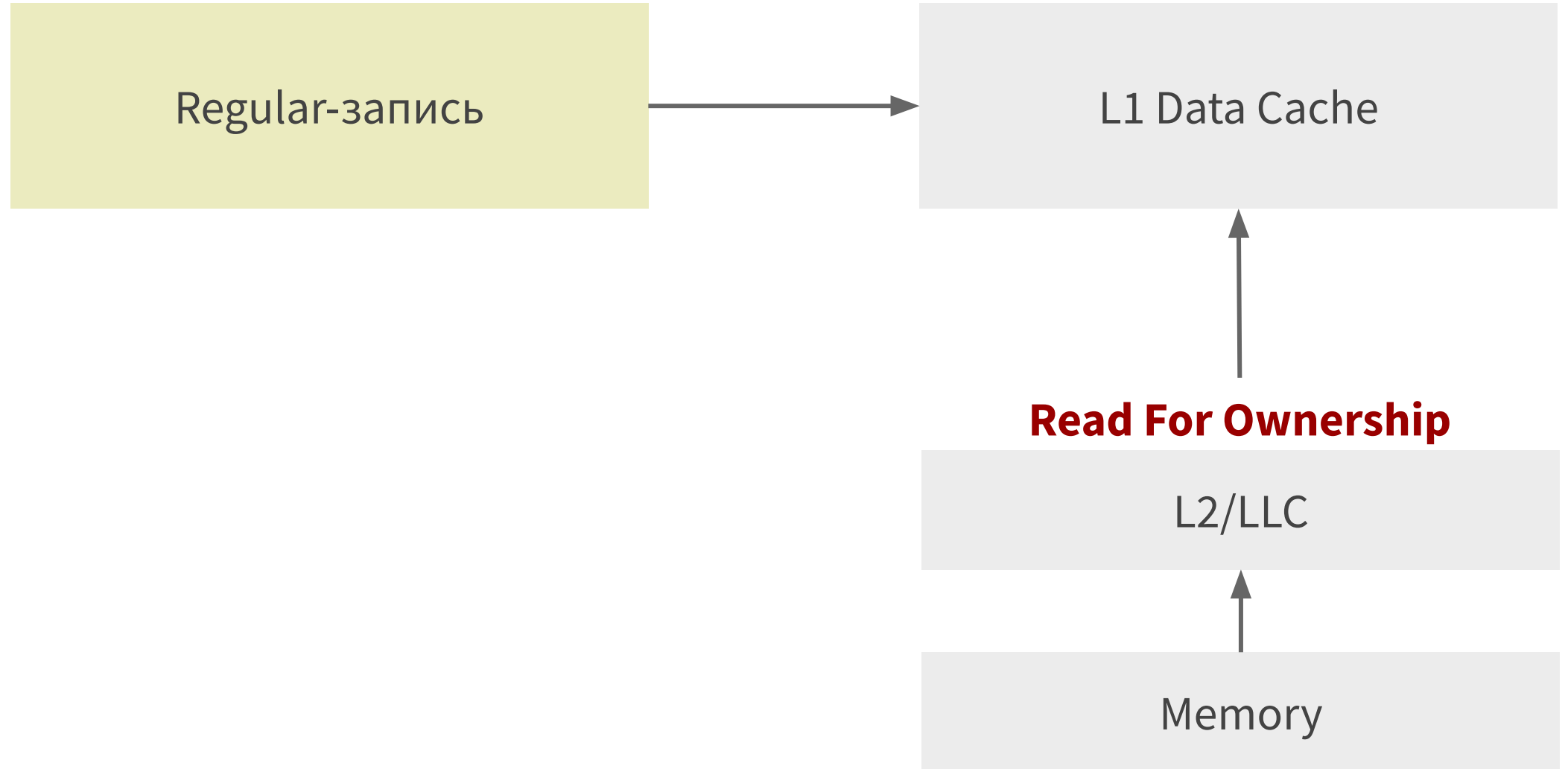
Write Back память: regular stores



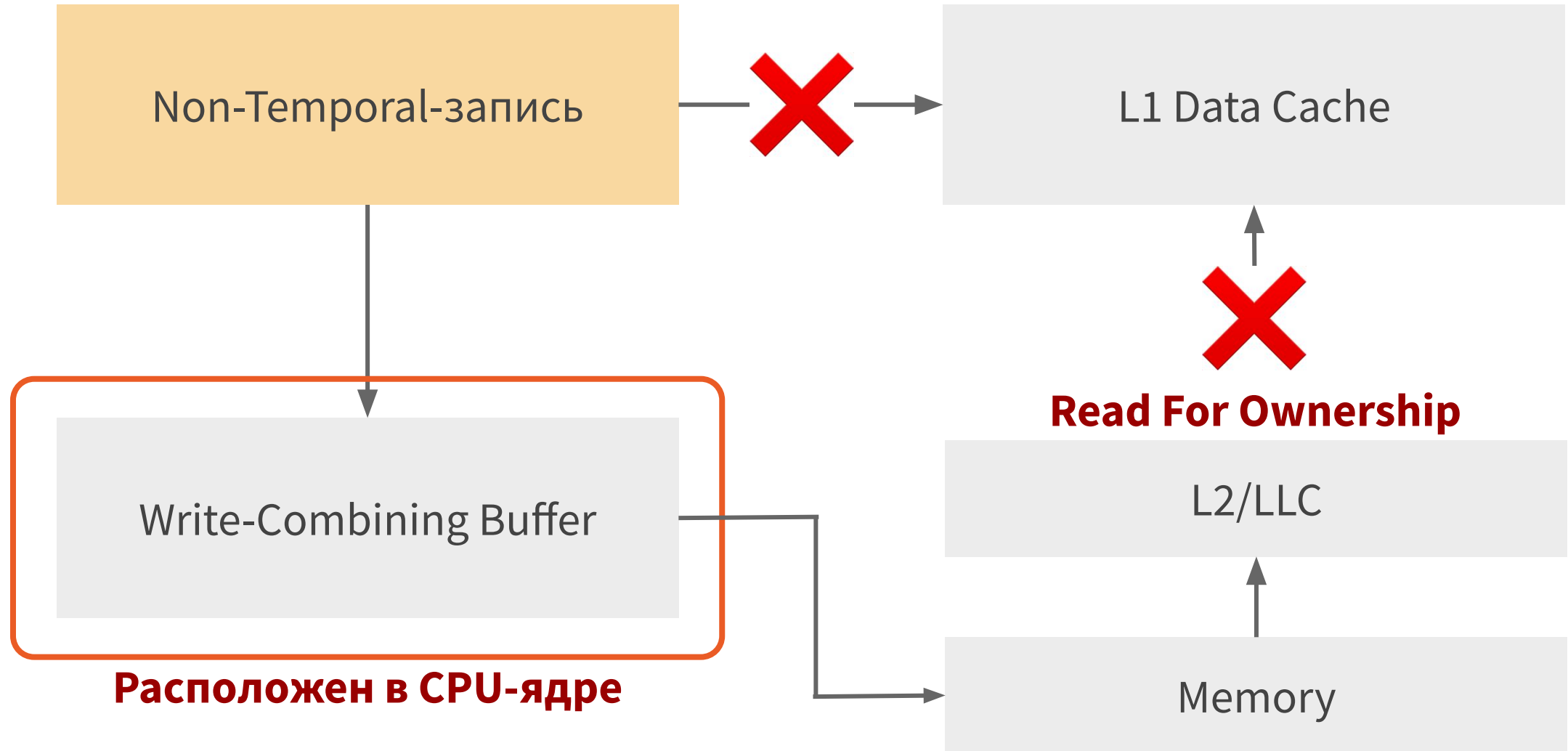
Write Back память: regular stores



Write Back память: regular stores



Write Back память: non-temporal stores



Libc memsru: техника копирования

- 4KB до 6 MB

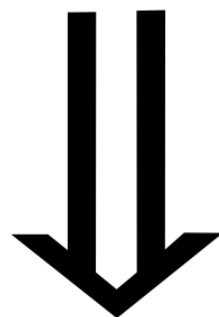
rep movsb

- Более 6MB

Non-Temporal stores

Non-temporal stores: перфоманс копирования

Нет RFO-трафика

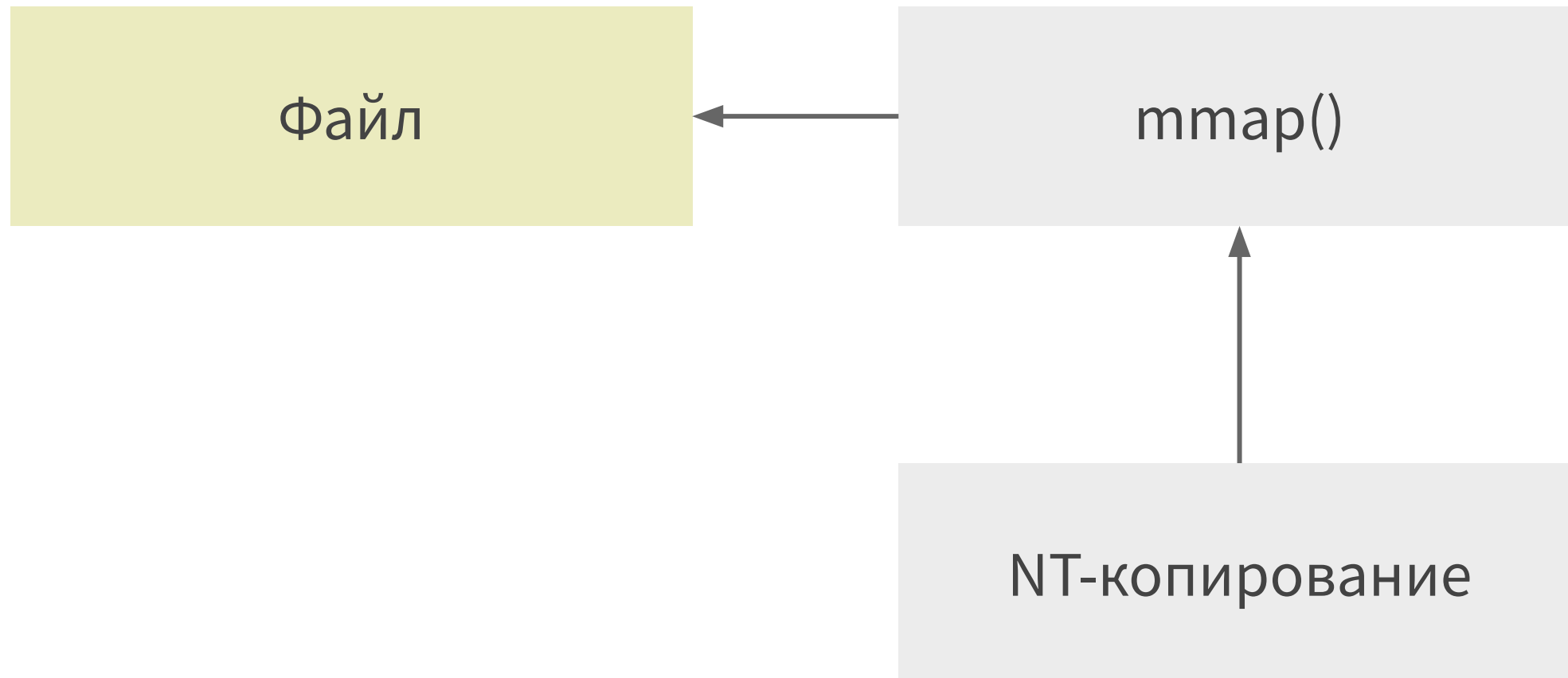


**Перфоманс копирования
больших данных**

Копирование памяти в ядре

- Почему не используется AVX?
- Будет ли AVX эффективнее `rep movsb`?
- **Если AVX быстрее, то можно ли его использовать?**

Чтение из файла: Non-Temporal stores + mmap



mpar: библиотека one-nio

```
1  String FILE_NAME;  
2  long mappedRegionSize;  
3  
4  one.nio.mem.MappedFile mf = new one.nio.mem.MappedFile(  
5      FILE_NAME,  
6      mappedRegionSize,  
7      one.nio.mem.MappedFile.MAP_RO  
8  );  
9  
10 mf.close();
```

mmap: библиотека one-nio

```
1 String FILE_NAME;  
2 long mappedRegionSize;  
3  
4 one.nio.mem.MappedFile mf = new one.nio.mem.MappedFile(  
5     FILE_NAME,  
6     mappedRegionSize,  
7     one.nio.mem.MappedFile.MAP_RO  
8 );  
9  
10 mf.close();
```

mmap()

mmap: библиотека one-nio

```
1  String FILE_NAME;  
2  long mappedRegionSize;  
3  
4  one.nio.mem.MappedFile mf = new one.nio.mem.MappedFile(  
5      FILE_NAME,  
6      mappedRegionSize,  
7      one.nio.mem.MappedFile.MAP_RO  
8  );  
9  
10 mf.close();
```



Non-Temporal Stores + mmap: перфоманс

- горячий page cache

Non-Temporal Stores + mmap: перфоманс

- горячий page cache

Без учета mmap/munmap

5.5GB	mmap+NT, sec	criticalRead, sec
	0.72	0.72

Non-Temporal Stores + mmap: перфоманс

- горячий page cache

Без учета mmap/munmap

	mmap+NT, sec	criticalRead, sec
5.5GB	0.72	0.72

С учетом mmap/munmap

	mmap+NT, sec	criticalRead, sec
5.5GB	0.82	0.72

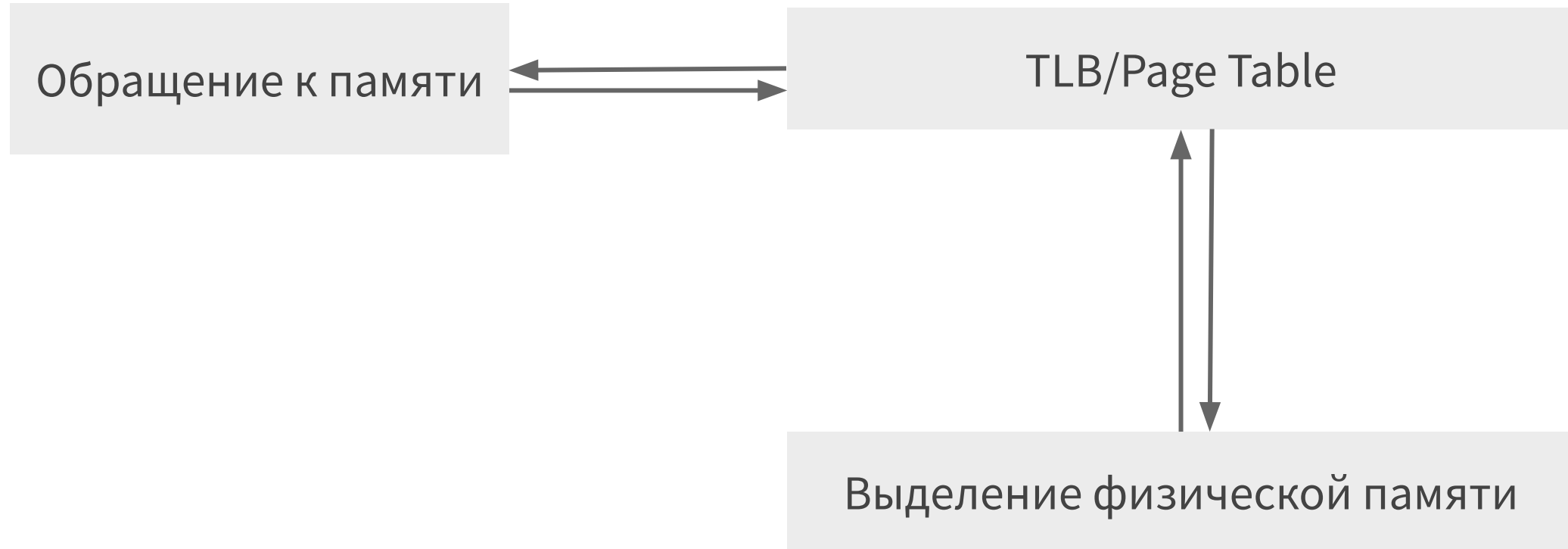
Non-Temporal Stores + mmap: анализ

```
$ sudo perf record --call-graph dwarf -F 9123 -p <app_pid>
```

```
- 29,50% __memmove_avx_unaligned_erms
  - 22,56% page_fault
    - 22,51% do_page_fault
      + 22,46% __do_page_fault
    + 3,84% swapgs_restore_regs_and_return_to_usermode
  + 1,79% error_entry
2,64% page_fault
```

Page Fault при первом доступе

Page Fault



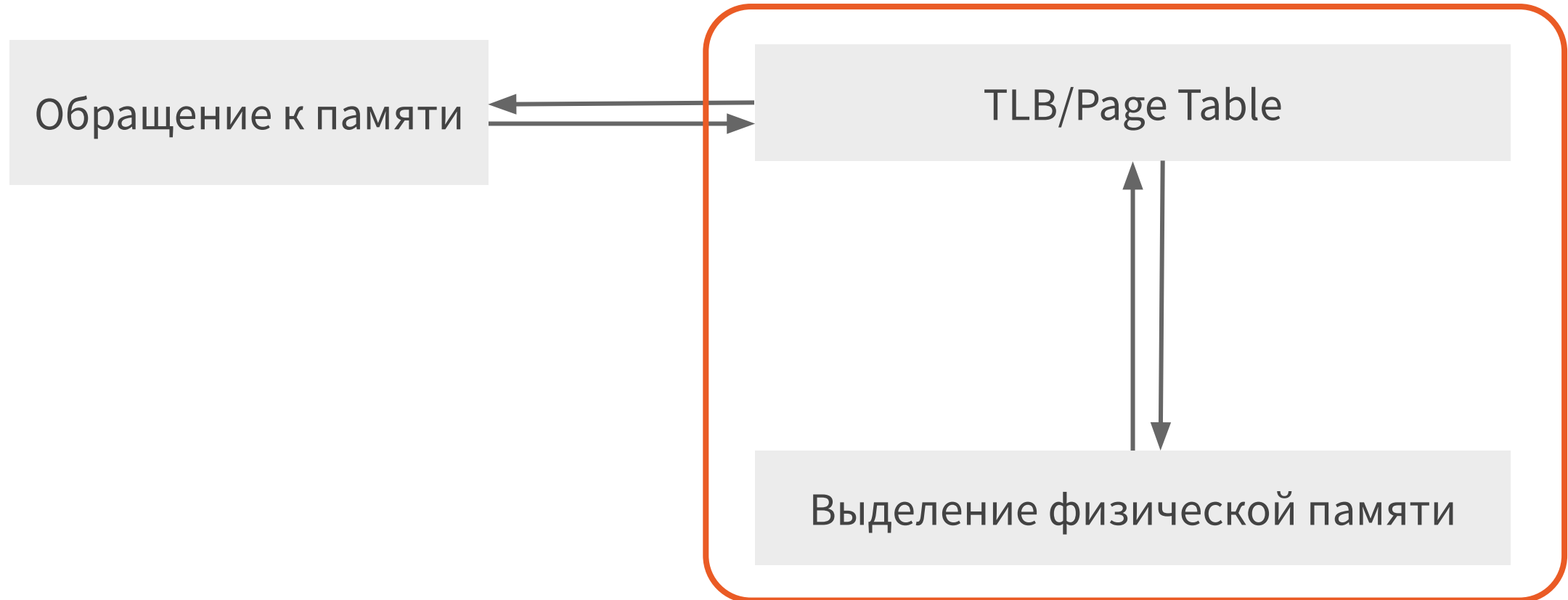
Page Fault

Трансляция виртуального адреса на физический



Page Fault

Minor Page Fault



Non-Temporal Stores + mmap: перфоманс

- горячий page cache

Без учета mmap/munmap

	mmap+NT, sec	criticalRead, sec
5.5GB	0.72	0.72

С учетом mmap/munmap

	mmap+NT, sec	criticalRead, sec
5.5GB	0.82	0.72

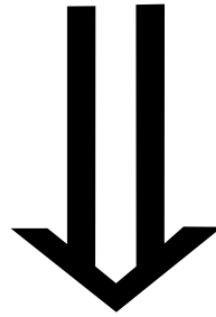
Без учета Page Fault

	mmap+NT, sec	criticalRead, sec
5.5GB	0.51	0.72

Non-Temporal Stores + mmap: резюме

Page Fault

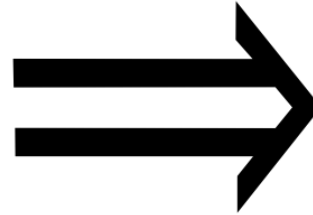
Cache Miss



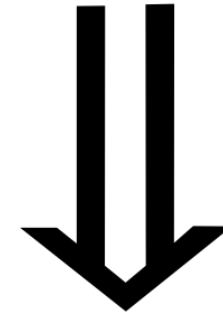
Не подходит для чтения
файла

Non-Temporal Stores: резюме

Не кэшируются



Уменьшение
загрязнения кэшей



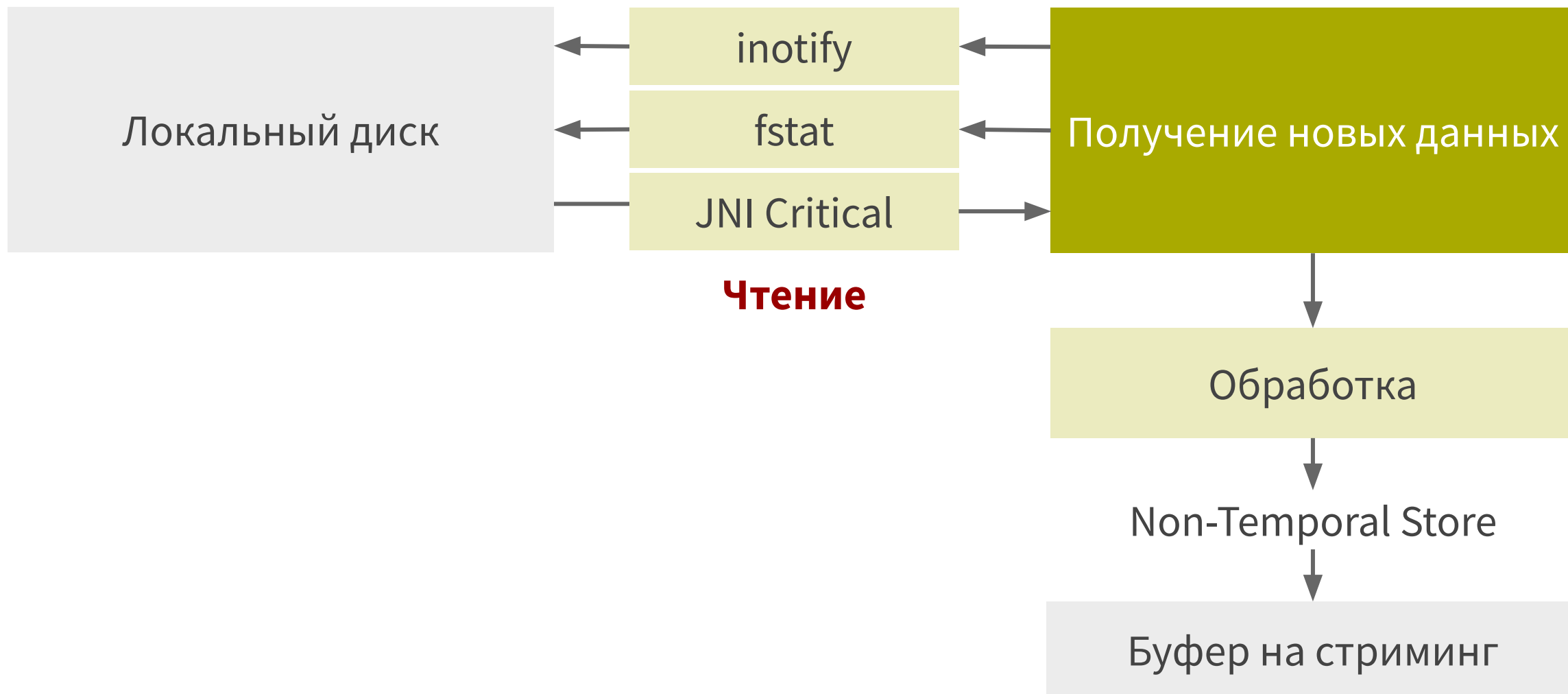
**Идеальный выбор
для стриминга**

Оптимизация производительности: резюме

- Определение изменившихся файлов ✓
- Определение предыдущей позиции стриминга ✓
- Чтение новых данных ✓
- Non-Temporal-стриминг обработанных данных ✓

Итоговая схема система стриминга

Мониторинг



Итоги

- Экономия на инфраструктуре до 40%
- Авто-скейлинг
- Уменьшены затраты на стриминг
 - Исключены накладные расходы на I/O
 - Оптимизировано использование кэша

Заключение

Повышение производительности за счет:

- Выхода за границы абстракций JDK
- Использования возможностей железа и O/S
- Обхода особенностей JVM

Q&A



Thank you!

www.griddynamics.com